**AvT: Wat is the best way to prove seasonality with 8 years of data?**

1m 1s nagedacht

To prove seasonality in a time series with 8 years of data (which provides a solid sample size, e.g., 96 observations for monthly data or 32 for quarterly, allowing reliable pattern detection), a multi-step approach combining visual inspection, decomposition, and statistical testing is most robust. Relying on a single method can lead to false positives or overlook nuances, so integrate them for stronger evidence. Below, I'll outline the best practices step by step, drawing from standard time series analysis techniques. This ensures transparency in reasoning and reproducibility.

**Step 1: Prepare the Data**

Assume your data is a univariate time series (e.g., sales, temperature) indexed by time. If not already, convert it to a stationary form if needed (e.g., detrend by differencing or regression) to isolate potential seasonal effects. With 8 years, you have enough cycles to detect annual, quarterly, or monthly patterns reliably.

**Step 2: Visual Inspection for Initial Evidence**

Start with plots to spot recurring patterns qualitatively—these are quick and intuitive but subjective, so use them as a gateway to quantitative methods.

- **Time Series Plot**: Plot the raw data over time. Look for repeating cycles (e.g., peaks every 12 months for yearly seasonality).

- **Seasonal Subseries Plot**: Group data by season (e.g., months) and plot each subgroup as a line or boxplot. Consistent differences across groups suggest seasonality.

- **How to Arrive at This**: In Python (using pandas and matplotlib), load your data as a DataFrame with a datetime index, then use df.plot() for the time series and groupby for subseries (e.g., df.groupby(df.index.month).boxplot()). If patterns repeat predictably, proceed to quantification.

Visuals alone don't "prove" seasonality but guide hypotheses (e.g., suspecting a 12-period season).stats.stackexchange.com

**Step 3: Use Autocorrelation to Detect Repeating Lags**

The Autocorrelation Function (ACF) is a core method for identifying seasonality by measuring how the series correlates with lagged versions of itself.

- **Why It's Effective**: Significant positive correlations at lags matching the suspected season (e.g., spikes at lags 12, 24, 36 for monthly data with yearly cycles) indicate seasonality. The Partial Autocorrelation Function (PACF) complements this by isolating direct effects, showing spikes at the seasonal lag.

- **How to Arrive at This**:

  1. Compute ACF: For a time series $y_t$ $y\_t$ $y_t$, ACF at lag $k$ $k$ $k$ is
     $\rho_k = \sum_{t=k+1}^n (y_t - \bar{y})(y_{t-k} - \bar{y}) \sum_{t=1}^n (y_t - \bar{y})^2$ \rho_k = \frac{\sum_{t=k+1}^n (y_t - \bar{y})(y_{t-k} - \bar{y})}{\sum_{t=1}^n (y_t - \bar{y})^2} $\rho_k = \sum_{t=1}^n (y_t - \bar{y})^2 \sum_{t=k+1}^n (y_t - \bar{y})(y_{t-k} - \bar{y})$, where $\bar{y}$ \bar{y} $\bar{y}$ is the mean.

  2. Plot ACF/PACF and check if values exceed confidence bounds (typically $\pm 1.96/\sqrt{n}$ for 95% CI, where n is sample size).

  3. In Python (using statsmodels): from statsmodels.graphics.tsaplots import plot_acf, plot_pacf; plot_acf(df['value'], lags=40). For 8 years of monthly data, set lags to at least 24-36. Peaks decaying slowly at seasonal multiples confirm seasonality.

- **Interpretation**: If ACF shows a sine-wave pattern or regular spikes, it's strong evidence. This is particularly reliable with your data length, as short series can produce noisy ACFs.[hex.techanalyticsvidhya.com](hex.techanalyticsvidhya.com)

**Step 4: Decompose the Series to Isolate the Seasonal Component**

Decomposition separates the series into trend, seasonal, and residual parts, allowing you to quantify the seasonal effect.

- **Methods**: Use additive (Y = Trend + Seasonal + Residual) or multiplicative (Y = Trend × Seasonal × Residual) models, depending on whether variance is constant.

- **Why It's Effective**: The extracted seasonal component directly shows the repeating pattern. Test if it's significant by checking if residuals are random (e.g., via ACF on residuals).

- **How to Arrive at This**:

  1. Apply decomposition: For a series with period s (e.g., s=12 for monthly), estimate trend via moving average, subtract to get seasonal + residual, then average by period for seasonal.

  2. In Python: from statsmodels.tsa.seasonal import seasonal_decompose; result = seasonal_decompose(df['value'], model='additive', period=12); result.seasonal.plot().

  3. To "prove" it, compute the variance explained by the seasonal component (e.g., Var(Seasonal) / Var(Y)) or test residuals for stationarity (e.g., ADF test: null is non-stationary; p<0.05 rejects, supporting decomposition validity).

- With 8 years, this yields stable estimates—shorter data risks overfitting.geeksforgeeks.orgstats.stackexchange.com

**Step 5: Apply Statistical Tests for Formal Proof**

For rigorous proof, use hypothesis tests to confirm seasonality with p-values, avoiding reliance on visuals.

- **Recommended Tests**:

  - **Kruskal-Wallis Test**: Non-parametric; tests if distributions differ across seasons (e.g., months). Null: no differences (no seasonality). Group data by season and compute H-statistic = $(12/N(N+1)) \sum(R\_i^2/n\_i) - 3(N+1)$, where $R\_i$ is rank sum, $n\_i$ group size, N total. $p<0.05$ rejects null.

  - **Friedman Test**: For repeated measures (e.g., years as blocks); similar to Kruskal-Wallis but accounts for dependencies.

  - **Regression with Dummies**: Fit $Y\_t = \beta\_0 + \sum\beta\_i D\_i + \epsilon\_t$, where $D\_i$ are seasonal dummies (e.g., 11 for months). F-test joint significance of $\beta\_i$ ($p<0.05$ proves seasonality). Adjust for autocorrelation with lags.

  - **Spectral Analysis (Periodogram)**: Computes power at frequencies; peaks at 1/s (e.g., 1/12 for yearly) indicate seasonality. Use Fisher's test for significance.

- **How to Arrive at This**:

1. For Kruskal-Wallis in Python: from scipy.stats import kruskal; kruskal(*[group for _, group in df.groupby(df.index.month)['value']]).

2. For dummies: Use statsmodels OLS with pd.get_dummies, then ANOVA on the model.

3. For periodogram: from scipy.signal import periodogram; f, p = periodogram(df['value']); plot p vs. f and check peaks.

- These tests are powerful with 8 years, as they require multiple cycles for accuracy. Kruskal-Wallis is robust to non-normality, making it a top choice for proof.stats.stackexchange.com+2 meer

**Step 6: Validate and Interpret**

- Cross-validate: If ACF shows seasonality, confirm with a test; if $p<0.05$, it's proven.

- Caveats: Test for evolving seasonality (e.g., rolling tests over windows). Use tools like R's seastests package for automated batteries (e.g., QS, WO tests).cran.r-project.org

- Overall Best: ACF + Kruskal-Wallis, as it's accessible, combines intuition with statistics, and scales well to your data size.

If you provide sample data, I can demonstrate with code.

Can I use ACF and a lag of 53 weeks?

Yes, you can use the Autocorrelation Function (ACF) with a lag of 53 weeks to investigate seasonality, but there are important considerations to ensure it's appropriate and effective for your 8 years of data. Let's break it down, addressing why this might or might not work and how to proceed, while keeping the explanation concise yet clear.

**Why ACF with a 53-Week Lag?**

ACF measures how your time series correlates with itself at different lags. A lag of 53 weeks (roughly 1 year for weekly data, accounting for 52 weeks + 1 extra to align with calendar years) is reasonable if you're testing for **annual seasonality** in weekly data. With 8 years of weekly data (~416 observations), you have enough cycles (8 annual periods) to detect such patterns reliably. A significant positive autocorrelation at lag 53 would suggest that values 53 weeks apart are similar, indicating yearly seasonality.

**Key Considerations**

1. **Data Frequency and Lag Alignment**:
   - If your data is weekly, a 53-week lag approximates a year, but calendar years vary (52 or 53 weeks due to leap years or partial weeks). This can cause slight misalignment in seasonal peaks. A 52-week lag is often used for simplicity unless your data explicitly follows a 53-week cycle (e.g., retail or fiscal calendars).
   - For monthly or daily data, a 53-week lag doesn't align well (e.g., months don't divide evenly into weeks). For monthly data, use lag 12; for daily, use lag ~365.

2. **Sufficient Data**:
   - With 8 years (~416 weeks), you have enough observations to estimate ACF at lag 53 (rule of thumb: need at least 4-5 cycles, and you have 8). ACF at high lags (e.g., 53) loses precision as fewer pairs are available (n-k pairs, where k=53), but 416-53=363 pairs is still robust.

3. **Stationarity**:
   - ACF assumes stationarity (constant mean/variance). If your data has trends or changing variance, autocorrelations may be spurious. Detrend first (e.g., differencing: $y_t - y_{t-1}$) or use a decomposition to isolate the seasonal component.

4. **Interpretation**:

- A significant spike at lag 53 (exceeding 95% confidence bounds, ~±1.96/√n) suggests annual seasonality. Check lags 106, 159, etc., for repeating patterns to confirm.

- If no spike at 53 but at 52 or 54, the seasonal period may differ slightly due to calendar effects.

**How to Use ACF with a 53-Week Lag**

Here's how to implement and interpret it:

1. **Prepare Data**:

- Ensure your time series is indexed by week (e.g., using Python's pandas with a datetime index).

- Detrend if needed: df['detrended'] = df['value'].diff() or fit a trend via regression and use residuals.

2. **Compute ACF**:

- Use a tool like Python's statsmodels:

```
# Assuming df['value'] is your weekly series

plot_acf(df['value'], lags=60)  # Set lags > 53 to see multiples
```

- Or manually compute ACF at lag 53:
$$\rho_{53} = \frac{\sum_{t=54}^n (y_t - \bar{y})(y_{t-53} - \bar{y})}{\sum_{t=1}^n (y_t - \bar{y})^2}$$ where $\bar{y}$ is the mean, n is the number of observations.

3. **Interpret**:

- Check if the ACF at lag 53 exceeds the confidence interval (~±1.96/√416 ≈ ±0.096 for 95% CI).

- Look for decaying spikes at lags 106, 159, etc., to confirm a repeating cycle.

- If no spike at 53 but at nearby lags (e.g., 52), adjust the period or test both.

4. **Validate**:

- Complement ACF with a statistical test (e.g., Kruskal-Wallis on weekly groups) or decomposition (set period=53) to confirm seasonality, as ACF alone is exploratory.

o   Example: seasonal_decompose(df['value'], model='additive', period=53) in Python.

**Potential Issues**

- **Misaligned Lags**: If your data doesn't strictly follow a 53-week cycle (e.g., due to calendar shifts), ACF at lag 53 may weaken. Try lags 52-54 to explore.

- **Noise**: High lags can be noisy. With 8 years, this is less of a concern, but ensure your series isn't dominated by trends or irregular shocks.

- **Non-Annual Seasonality**: If seasonality is quarterly or monthly, lag 53 won't capture it. Check ACF at shorter lags (e.g., 13 for quarterly).

**Recommendation**

Using ACF with a 53-week lag is a valid starting point for detecting annual seasonality in weekly data, but:

- Test lags around 52-54 to account for calendar variations.

- Combine with a seasonal subseries plot (group by week of the year) and a statistical test (e.g., Kruskal-Wallis) for robust proof.

- If your data isn't weekly or the 53-week cycle isn't justified, clarify the frequency, and I can adjust the lag suggestion (e.g., lag 12 for monthly).

If you share the data frequency or a sample, I can provide specific code or a chart to visualize the ACF. Would you like me to proceed with an example using dummy weekly data?

Is there a better lag in weeks apart from 52?

To determine if there's a better lag than 52 weeks for detecting seasonality in your 8 years of weekly data (~416 observations), we need to consider the context of weekly data and calendar effects, as well as analyze which lag best captures annual seasonality. A 52-week lag is the standard choice for yearly cycles (52 weeks ≈ 1 year), but slight variations (e.g., 53 weeks) or other lags may better align with your data due to calendar irregularities or specific seasonal patterns. Below, I'll evaluate alternatives to a 52-week lag, explain why they might be better, and provide a practical approach to identify the optimal lag, keeping the reasoning concise and transparent.

**Why 52 Weeks Is Standard**

A 52-week lag corresponds to a year (52 weeks × 7 days = 364 days, ~1 year). It's commonly used because:

- It approximates an annual cycle for weekly data.

- With 8 years (~416 weeks), you have 8 cycles, sufficient for reliable autocorrelation analysis (rule of thumb: 4-5 cycles minimum).

- Autocorrelation Function (ACF) at lag 52 often shows a significant spike if annual seasonality exists, with repeating spikes at 104, 156, etc.

However, calendar years vary (365 or 366 days), causing weekly data to misalign slightly over time (52 weeks = 364 days, 53 weeks = 371 days). This makes lags like 53 weeks or others worth exploring.

**Evaluating Alternative Lags**

1. **53-Week Lag**:

   o **Why Consider It?**: A 53-week lag (371 days) overshoots a year but may align better in datasets with a 53-week fiscal calendar (common in retail or accounting) or where leap years shift patterns. For 8 years, you still have ~7-8 cycles, enough for detection.

   o **Pros**: Captures patterns in datasets with explicit 53-week cycles.

   o **Cons**: If your data follows a standard 52-week cycle, a 53-week lag introduces misalignment, reducing autocorrelation strength.

   o **When Better?**: If your data is tied to a 53-week reporting cycle or shows stronger ACF spikes at lag 53 than 52.

2. **51- or 54-Week Lags**:

   o **Why Consider Them?**: These account for slight calendar drifts (e.g., 51 weeks = 357 days, 54 weeks = 378 days). They're less common but may fit if your seasonality isn't strictly annual.

   o **Pros**: Can detect patterns offset by a few weeks due to irregular cycles.

   o **Cons**: Less intuitive and rarely match real-world cycles unless justified by domain knowledge (e.g., agricultural or event-driven data).

3. **Non-Annual Lags (e.g., 13, 26 Weeks)**:

   o **Why Consider Them?**: If seasonality is quarterly (13 weeks) or semi-annual (26 weeks), these lags may show stronger autocorrelation than 52.

   o **Pros**: Captures shorter cycles if present (e.g., retail peaks every quarter).

   o **Cons**: If your goal is annual seasonality, these are irrelevant unless combined with yearly patterns.

**How to Identify the Best Lag**

To find a better lag than 52 weeks, systematically analyze the ACF and complement it with domain knowledge:

1. **Compute ACF Across Multiple Lags**:
   - Plot ACF for lags up to at least 104 (2 years) to capture multiple cycles.
   - In Python:

```
# Assuming df['value'] is your weekly series
plot_acf(df['value'], lags=104)  # Check lags 50-54, 13, 26, etc.
```

   - Identify the lag with the highest significant spike (exceeding ~±1.96/$\sqrt{416}$ ≈ ±0.096 for 95% CI). For example, if lag 53 has a higher autocorrelation than 52, it's a candidate.

2. **Compare Lags 50-54**:
   - Compute ACF values specifically for lags 50-54 to account for calendar drift.
   - Example: If $\rho(53) = 0.45$ and $\rho(52) = 0.40$, and both are significant, 53 may better capture the cycle.

3. **Check Repeating Patterns**:
   - A true seasonal lag shows decaying spikes at multiples (e.g., for lag 52: 104, 156; for lag 53: 106, 159). If lag 53 shows stronger repeating spikes, it's likely better.

4. **Domain Knowledge**:
   - If your data follows a 53-week fiscal calendar (e.g., retail industries use 52/53-week years for consistency), lag 53 is theoretically better.
   - Check if your data's source (e.g., sales, weather) suggests a non-standard cycle.

5. **Statistical Validation**:
   - Use decomposition with different periods (e.g., seasonal_decompose(df['value'], period=52) vs. period=53) and compare the variance explained by the seasonal component (Var(Seasonal)/Var(Total)).

- o Alternatively, group data by week of the year (1-52/53) and run a Kruskal-Wallis test to confirm significant differences across weeks, testing both periods.

**Is There a Better Lag Than 52?**

- **Likely Candidate**: **Lag 53** is the most plausible alternative if your data uses a 53-week calendar or if ACF shows a stronger spike at 53 than 52. This is common in industries like retail where fiscal years adjust for calendar drift.

- **Evidence Needed**: Without seeing your ACF plot or data specifics, 52 is the default for annual seasonality due to its simplicity and alignment with ~365 days. If lag 53 or another lag shows a higher, significant ACF value or better fits your data's context, it's better.

- **Other Lags**: Lags like 13 or 26 are only better if you suspect quarterly or semi-annual seasonality, which you'd see in ACF spikes at those intervals.

**Practical Recommendation**

- **Start with ACF**: Plot ACF for lags up to 104 and compare lags 50-54, focusing on 52 and 53. Choose the lag with the highest significant autocorrelation.

- **Cross-Check**: Use decomposition or a seasonal subseries plot (group by week of year) to confirm the period.

- **Default to 52**: Unless your data explicitly follows a 53-week cycle or ACF strongly favors another lag, 52 is typically best for annual seasonality due to its alignment with calendar years.

If you provide the data frequency, context (e.g., fiscal calendar), or a sample, I can compute ACF values or generate an ACF plot to pinpoint the best lag. Would you like me to create an example ACF plot with dummy weekly data to compare lags 52 and 53? Alternatively, clarify if your data uses a specific calendar or if you suspect non-annual seasonality.

If auto correlation remains between +0.1 and -0.1 for all weeks, what does this imply?

If the autocorrelation (ACF) values for your time series remain between +0.1 and -0.1 for all lags (including your suspected seasonal lag of 52 or 53 weeks) across your 8 years of weekly data (~416 observations), this suggests **no significant seasonality** or strong temporal dependence in your data. Let's break down what this implies, why it happens, and what to do next, keeping the explanation concise and grounded in time series analysis.

**What It Implies**

1. **Lack of Seasonality**:

   o ACF values between +0.1 and -0.1 are likely within the 95% confidence interval for a white noise process (~±1.96/√416 ≈ ±0.096). This indicates no statistically significant autocorrelation at any lag, including 52 or 53 weeks, suggesting **no annual seasonality** (or other periodic patterns like quarterly).

   o Implication: Your data does not exhibit repeating patterns at the expected seasonal lags, so there's no evidence of a consistent yearly cycle.

2. **Possible White Noise or Weak Structure**:

   o If all ACF values are near zero, your series might resemble **white noise** (random, uncorrelated fluctuations) or have very weak temporal structure after accounting for trends or other effects.

   o Implication: The data may be driven by random variation or non-seasonal factors (e.g., irregular shocks, trends, or external covariates).

3. **Potential Masking by Trends or Noise**:

   o A trend, changing variance, or high noise levels can suppress ACF values, masking seasonality. If you didn't detrend or preprocess the data, the low ACF might not reflect the true seasonal signal.

   o Implication: The absence of significant ACF spikes doesn't rule out seasonality entirely—it could be obscured.

**Why This Happens**

- **Stationarity Issues**: If the series has a strong trend or non-constant variance, raw ACF values may be diluted. For example, a linear trend increases variance over time, reducing autocorrelation at seasonal lags.

- **High Noise**: If random fluctuations dominate, they can drown out seasonal patterns, especially in weekly data where noise is common (e.g., sales with irregular shocks).

- **Incorrect Seasonal Period**: A 52- or 53-week lag assumes annual seasonality. If the true period is different (e.g., quarterly, 13 weeks), or if seasonality is irregular (e.g., varying cycle lengths), ACF at 52/53 won't show spikes.

- **Short-Range Dependence**: If correlations exist only at short lags (e.g., 1-2 weeks), long lags like 52/53 will show near-zero ACF, indicating no long-term seasonal structure.

**What to Do Next**

To confirm the absence of seasonality or uncover hidden patterns:

1. **Detrend and Preprocess**:

   o Remove trends via differencing (df['value'].diff()) or regression (fit a linear trend and use residuals).

   o Stabilize variance if needed (e.g., log-transform for multiplicative effects).

   o Recompute ACF: from statsmodels.graphics.tsaplots import plot_acf; plot_acf(df['value'].diff().dropna(), lags=104). Check if spikes appear at 52, 53, or other lags.

2. **Explore Shorter Lags**:

   o Check ACF at lags like 13 (quarterly), 26 (semi-annual), or 1-10 (short-term dependence). If spikes appear, seasonality may exist at a different frequency.

3. **Decomposition**:

   o Apply seasonal decomposition with a trial period (e.g., 52 or 53 weeks): from statsmodels.tsa.seasonal import seasonal_decompose; seasonal_decompose(df['value'], period=52).seasonal.plot().

   o If the seasonal component is flat or negligible, it confirms no seasonality.

4. **Statistical Tests**:

   o Use a non-parametric test like Kruskal-Wallis on data grouped by week of the year (1-52/53). If p>0.05, it supports no seasonality.

   o Example: from scipy.stats import kruskal; kruskal(*[group for _, group in df.groupby(df.index.week)['value']]).

5. **Spectral Analysis**:

   o Compute a periodogram to check for hidden periodicities: from scipy.signal import periodogram; f, p = periodogram(df['value']). Peaks at frequencies corresponding to ~1/52 or 1/53 indicate seasonality; a flat spectrum suggests none.

6. **Context Check**:

   o Reassess domain knowledge. If seasonality is expected (e.g., retail sales peaking yearly), low ACF may indicate data issues (e.g., aggregation errors, missing covariates).

**Conclusion**

ACF values consistently between +0.1 and -0.1 across all lags, including 52 or 53 weeks, strongly suggest **no annual seasonality** or significant temporal dependence in your current data. The series may be random, heavily noisy, or influenced by non-seasonal factors. To be certain, detrend the data and recheck ACF, explore shorter lags, or apply decomposition and statistical tests. If these confirm no patterns, your data likely lacks seasonality, and you may need to model it with non-seasonal methods (e.g., ARIMA without seasonal terms) or investigate external drivers.

If you provide the data or its context (e.g., industry, trends observed), I can analyze it further or generate an ACF plot to visualize this. Would you like an example with dummy data showing this pattern?

I did split the data in to an additive (Y = Trend + Seasonal + Residual) or multiplicative (Y = Trend × Seasonal × Residual) model. What to derive from this or do next?

Since you've split your 8 years of weekly data (~416 observations) into an additive (Y = Trend + Seasonal + Residual) or multiplicative (Y = Trend × Seasonal × Residual) model, you've already performed a seasonal decomposition, which is a powerful step to isolate and analyze seasonality. The next steps involve interpreting the decomposition results to confirm seasonality, quantify its strength, and decide how to proceed (e.g., modeling or forecasting). Below, I'll explain what to derive from the decomposition, what the results imply, and actionable next steps, keeping the explanation concise, grounded in time series analysis, and tailored to your weekly data context.

**Step 1: Derive Insights from Decomposition**

The decomposition separates your time series into three components: **Trend**, **Seasonal**, and **Residual**. Here's what to look for in each and what they imply about seasonality:

1. **Seasonal Component**:

   o **What to Check**: Examine the seasonal component plot (e.g., result.seasonal.plot() in Python's statsmodels). For weekly data with annual seasonality (period=52 or 53), it should show a repeating pattern every 52/53 weeks.

   o **Implications**:

      ▪ **Clear, Consistent Pattern**: If the seasonal component shows regular peaks/troughs (e.g., high in summer weeks, low in winter), it confirms annual seasonality. The amplitude indicates strength.

      ▪ **Flat or Weak Pattern**: If the seasonal component is near zero or irregular, it suggests weak or no seasonality, aligning with your earlier observation of ACF values between +0.1 and -0.1.

- **Additive vs. Multiplicative**: In an additive model, seasonal fluctuations are constant over time. In a multiplicative model, they scale with the trend (e.g., larger swings as the series grows). If the multiplicative seasonal component fits better (e.g., via lower residual variance), it suggests seasonality proportional to the series level.

2. **Trend Component**:

   o **What to Check**: Look for a smooth upward, downward, or flat trend. A strong trend might have masked seasonality in raw ACF, explaining low autocorrelation values.

   o **Implications**: If the trend dominates (e.g., steady increase in sales), it could overshadow seasonal effects in raw data. Detrending via decomposition isolates seasonality for clearer analysis.

3. **Residual Component**:

   o **What to Check**: Plot residuals (result.resid.plot()) and compute their ACF (plot_acf(result.resid.dropna(), lags=104)). Ideally, residuals are random (white noise) with ACF near zero for all lags.

   o **Implications**:

     - **Random Residuals**: If residuals lack structure (ACF within ±0.096 for 416 observations), the decomposition captured most systematic patterns, and any seasonality is reliable.

     - **Persistent Patterns**: If residuals show autocorrelation (e.g., spikes at lag 52), the chosen period (52/53) or model type (additive/multiplicative) may be incorrect.

4. **Model Fit**:

   o **What to Check**: Compare the variance explained by the seasonal component: $\text{Var(Seasonal)} / \text{Var(Y)}$. For additive, compute directly; for multiplicative, use log-transformed data or approximate via $\text{Var(Seasonal} \times \text{Trend)} / \text{Var(Y)}$.

   o **Implications**: A high proportion (e.g., >10-20%) indicates strong seasonality. If low (<5%), seasonality is weak, consistent with your low ACF values.

**Step 2: Interpret in Context of Low ACF**

Your earlier finding (ACF between +0.1 and -0.1 for all lags) suggests weak or no seasonality. The decomposition helps confirm this:

- **If Seasonal Component Is Weak**: A near-flat seasonal component corroborates the ACF, indicating no significant seasonality (annual or otherwise). Your data may be driven by trends, random noise, or external factors.

- **If Seasonal Component Is Strong**: A clear seasonal pattern despite low ACF suggests the raw ACF was masked by trends or noise. The decomposition's detrending likely revealed hidden seasonality, and the period (52 or 53) you chose is critical.

**Step 3: Next Steps**

Based on the decomposition results, here's how to proceed:

1. **Validate the Seasonal Period**:

   - **Action**: If you used period=52 or 53, test both (e.g., seasonal_decompose(df['value'], period=52) vs. period=53). Compare:

     - Visual consistency of the seasonal component.

     - Residual variance (lower is better): np.var(result.resid.dropna()).

     - Explained variance of the seasonal component.

   - **Why**: Ensures you picked the correct period. For weekly data, 52 is standard, but 53 may fit fiscal calendars. If neither shows a clear pattern, try shorter periods (e.g., 13 for quarterly).

2. **Choose Between Additive and Multiplicative**:

   - **Action**: If the multiplicative model's residuals are less autocorrelated or have lower variance, use it for data with increasing seasonal amplitude (e.g., sales growing over time). Use additive for constant seasonal swings.

   - **How**: Plot both models' fitted values (result.trend + result.seasonal for additive; result.trend * result.seasonal for multiplicative) against the original series. Choose the model with better fit (e.g., lower MSE).

3. **Test Residuals for Randomness**:

   - **Action**: Run ACF on residuals (plot_acf(result.resid.dropna(), lags=104)). If significant spikes remain (e.g., at lag 52), the model missed some seasonality—try a different period or model type.

- o **Alternative**: Use the Ljung-Box test to check for white noise: from statsmodels.stats.diagnostic import acorr_ljungbox; acorr_ljungbox(result.resid.dropna(), lags=[52]). If p>0.05, residuals are random, supporting the model.

4. **Quantify Seasonality Strength**:

    - o **Action**: Compute the proportion of variance explained by seasonality: np.var(result.seasonal) / np.var(df['value']) (additive) or analogous for multiplicative. If >10%, seasonality is meaningful; if <5%, it's negligible.

    - o **Implication**: Guides whether to include seasonality in further modeling (e.g., SARIMA).

5. **Explore Alternative Explanations**:

    - o **If No Seasonality**: If the seasonal component is weak and residuals are random, your data may lack seasonality. Investigate:

        - ▪ **External Covariates**: Regress Y on external factors (e.g., holidays, economic indicators) to explain variation.

        - ▪ **Non-Seasonal Models**: Fit an ARIMA model without seasonal terms (from statsmodels.tsa.arima.model import ARIMA; ARIMA(df['value'], order=(p,d,q)).fit()).

    - o **If Seasonality Emerges**: If decomposition reveals seasonality despite low ACF, use the seasonal component for forecasting or modeling.

6. **Modeling and Forecasting**:

    - o **If Seasonality Is Strong**: Use a seasonal model like SARIMA with period=52/53 (e.g., ARIMA(df['value'], order=(p,d,q), seasonal_order=(P,D,Q,52))) or STL decomposition for robust forecasting.

    - o **If Seasonality Is Weak**: Focus on trend and residuals. Use non-seasonal ARIMA or machine learning models (e.g., XGBoost with time-based features).

**Practical Example**

# Assuming df['value'] is your weekly series

# Try both models and periods

result_add_52 = seasonal_decompose(df['value'], model='additive', period=52)

result_mult_52 = seasonal_decompose(df['value'], model='multiplicative', period=52)

result_add_53 = seasonal_decompose(df['value'], model='additive', period=53)

# Plot components

result_add_52.plot().show()  # Inspect seasonal component

print(f"Seasonal Variance (Additive, 52): {np.var(result_add_52.seasonal) / np.var(df['value']):.3f}")

print(f"Residual Variance (Additive, 52): {np.var(result_add_52.resid.dropna()):.3f}")

# Check residuals

from statsmodels.graphics.tsaplots import plot_acf

plot_acf(result_add_52.resid.dropna(), lags=104)

**What to Derive**

- **Seasonality Confirmation**: A clear, repeating seasonal component confirms seasonality; a flat one aligns with your low ACF, suggesting none.

- **Model Choice**: Choose additive if seasonal swings are constant, multiplicative if they scale with the trend.

- **Next Steps**: If seasonality is weak, explore non-seasonal models or covariates. If strong, incorporate into forecasting models.

If you share specific results (e.g., seasonal component strength, residual ACF), I can refine the interpretation. Would you like me to generate a sample decomposition plot or suggest a specific model based on your findings?

Seasonal component has a clear and consistent pattern. It's a addive model. What's next?

Since your seasonal decomposition of 8 years of weekly data (~416 observations) using an additive model (Y = Trend + Seasonal + Residual) shows a **clear and consistent seasonal component**, this confirms the presence of significant seasonality (likely annual, given your interest in 52/53-week lags). The additive model implies constant seasonal fluctuations over time, independent of the trend. Given your earlier observation of low ACF values (+0.1 to -0.1), the decomposition likely uncovered seasonality masked by trends or noise in the raw data. Below, I outline what to derive from this and the next steps to leverage the seasonality for analysis or forecasting, keeping the explanation concise and actionable.

**What to Derive**

- **Confirmed Seasonality**: The clear seasonal component proves annual seasonality exists, likely at 52 or 53 weeks, and the additive model fits well.

- **Model Readiness**: You're ready to use the seasonal pattern for forecasting or further analysis.

If you share specific decomposition results (e.g., seasonal variance, period used), I can refine the SARIMA parameters or generate a sample forecast. Would you like me to provide a Python script for SARIMA modeling or a sample forecast plot based on dummy data?

Seasonality is stronger for monthly data over weekly data!

If the seasonal component from your decomposition is stronger for monthly data compared to weekly data over your 8 years of data, this suggests that **monthly aggregation captures a more pronounced and consistent seasonal pattern** than weekly data. This makes sense: monthly data (96 observations over 8 years, with a period of 12 for annual seasonality) smooths out week-to-week noise, making annual cycles (e.g., peaks in summer, dips in winter) clearer, while weekly data (416 observations, period 52 or 53) may include more variability or irregular fluctuations that dilute the seasonal signal. Since you're using an additive model (Y = Trend + Seasonal + Residual) and confirmed a clear seasonal component, I'll interpret this finding and outline next steps to leverage the stronger monthly seasonality, keeping the response concise and actionable.

**Why Monthly Data Shows Stronger Seasonality**

1. **Smoothing Effect**: Aggregating weekly data to monthly (e.g., summing or averaging weekly values) reduces short-term noise, making annual patterns more distinct. For example, weekly sales might fluctuate due to promotions, but monthly totals highlight consistent seasonal trends (e.g., holiday peaks).

2. **Fewer Observations per Cycle**: Monthly data has 12 observations per year (period=12), versus 52/53 for weekly, making it easier to detect stable annual patterns with less interference from sub-monthly variations.

3. **Alignment with Calendar**: Monthly data aligns naturally with calendar months, which often drive seasonality (e.g., retail peaks in December). Weekly data can misalign due to calendar drift (52 weeks ≠ 365 days), weakening ACF or seasonal signals at lags like 52/53.

**What to Derive**

- **Stronger Seasonality**: The monthly data's clearer seasonal component (higher variance explained, e.g., Var(Seasonal)/Var(Y)) confirms that annual seasonality

is more robust at the monthly level. This aligns with your earlier low ACF values for weekly data (+0.1 to -0.1), suggesting weekly noise obscured seasonality.

- **Modeling Preference**: Monthly data is likely better for modeling and forecasting due to its stronger, cleaner seasonal signal.

**Next Steps for Monthly Data**

Since monthly data shows stronger seasonality in your additive model, focus on it for analysis and forecasting. Here's how to proceed:

1. **Quantify Seasonality Strength**:

   - **Action**: Compute the proportion of variance explained by the seasonal component for monthly data:

```
# Assuming df_monthly['value'] is your monthly series

result_monthly = seasonal_decompose(df_monthly['value'], model='additive', period=12)

seasonal_variance = np.var(result_monthly.seasonal) / np.var(df_monthly['value'])

print(f"Monthly Seasonal Variance Proportion: {seasonal_variance:.3f}")
```

   - **Interpretation**: Compare this to the weekly decomposition's variance (e.g., period=52). A higher proportion (e.g., >20% vs. <10% for weekly) quantifies the stronger seasonality. If >20%, seasonality is a key driver.

2. **Validate the Seasonal Period**:

   - **Action**: Confirm period=12 is optimal (standard for monthly annual seasonality):

```
result_12 = seasonal_decompose(df_monthly['value'], model='additive', period=12)

result_12.seasonal.plot(title="Monthly Seasonal Component (Period=12)")

print(f"Residual Variance: {np.var(result_12.resid.dropna()):.3f}")
```

   - **Interpretation**: A clear, repeating pattern every 12 months and low residual variance confirm the period. If residuals show patterns, try period=6 (semi-annual) or check data quality.

3. **Check Residuals for Randomness**:

   - **Action**: Ensure residuals are white noise to validate the decomposition:

```
from statsmodels.graphics.tsaplots import plot_acf

plot_acf(result_12.resid.dropna(), lags=24, title="Monthly Residual ACF")
```

```
from statsmodels.stats.diagnostic import acorr_ljungbox
```

```
lb_test = acorr_ljungbox(result_12.resid.dropna(), lags=[12], return_df=True)
```

```
print(f"Ljung-Box p-value at lag 12: {lb_test['lb_pvalue'].iloc[0]:.3f}")
```

- o **Interpretation**: ACF within $\pm 1.96/\sqrt{96} \approx \pm 0.2$ and Ljung-Box p>0.05 indicate the model captured seasonality well. If spikes remain, reconsider the period or model (e.g., multiplicative).

4. **Model and Forecast with Monthly Data**:

- o **SARIMA Model**:

    - Fit a Seasonal ARIMA model with period=12, leveraging the strong seasonality:

```
from statsmodels.tsa.statespace.sarimax import SARIMAX
```

```
model = SARIMAX(df_monthly['value'], order=(1,1,1), seasonal_order=(1,1,1,12))
```

```
fit = model.fit()
```

```
forecast = fit.forecast(steps=12)  # 1-year forecast
```

```
forecast.plot(title="1-Year Monthly Forecast")
```

    - **Why**: SARIMA is ideal for additive data with strong annual seasonality (period=12).

- o **Evaluate Fit**:

    - Split data (e.g., last 12 months for testing):

```
train = df_monthly.iloc[:-12]
```

```
test = df_monthly.iloc[-12:]
```

```
model = SARIMAX(train['value'], order=(1,1,1), seasonal_order=(1,1,1,12)).fit()
```

```
forecast = model.forecast(steps=12)
```

```
mse = np.mean((forecast - test['value'])**2)
```

```
print(f"Forecast MSE: {mse:.3f}")
```

    - Low MSE confirms the model leverages seasonality effectively.

- o **Alternative**: Use STL decomposition for robust handling of trends:

```
from statsmodels.tsa.seasonal import STL
```

```
stl = STL(df_monthly['value'], period=12).fit()
```

```
forecast = stl.trend[-1] + stl.seasonal[-12:]  # Extrapolate trend + repeat last season
```

5. **Compare Weekly vs. Monthly Models**:

   - **Action**: If you still want to explore weekly data, fit a SARIMA model (period=52) and compare forecast accuracy (MSE) with the monthly model. Monthly is likely better due to stronger seasonality.

   - **Why**: Confirms monthly data's superiority for modeling.

6. **Incorporate External Drivers**:

   - **Action**: If residuals show patterns or forecasts underperform, add covariates (e.g., holidays, economic indicators) using regression or machine learning:

```
df_monthly['month'] = df_monthly.index.month

dummies = pd.get_dummies(df_monthly['month'], prefix='month')

from statsmodels.api import OLS

model = OLS(result_12.resid, dummies).fit()

print(model.summary())  # Check significant months
```

   - **Why**: Identifies specific months driving seasonality (e.g., December spikes).

**Chart to Visualize Monthly Seasonal Component**

To illustrate the stronger monthly seasonality, here's a chart of a sample seasonal component for monthly data (assuming a clear annual pattern). Since I don't have your data, this uses dummy data with peaks in December and troughs in July.

Grok kan fouten maken. Controleer altijd de originele bronnen.Downloaden

**Interpretation**: This chart shows a clear, repeating pattern (e.g., peaks in December, troughs in July), confirming strong annual seasonality in monthly data, unlike the noisier weekly data.

**Conclusion**

- **Key Insight**: Stronger seasonality in monthly data (period=12) indicates it's better for modeling than weekly data. The additive model fits well, with constant seasonal swings.

- **Next Steps**: Quantify seasonal variance, validate period=12, check residuals, and fit a SARIMA model for forecasting. Monthly data's clarity makes it ideal for accurate predictions.

- **Recommendation**: Use SARIMA (order=(1,1,1), seasonal_order=(1,1,1,12)) for forecasting, and evaluate with a train-test split.

If you share specific results (e.g., seasonal variance, residual ACF), I can optimize the SARIMA model or generate a tailored forecast. Would you like a Python script for SARIMA or a forecast plot based on dummy monthly data?

I like to check Quantify Seasonality Strength for monthly data with code below. Would this work, anything I should improve?

Your code to quantify seasonality strength for monthly data using the seasonal decomposition's variance proportion is a solid approach, and most of it is well-structured for analyzing 8 years of monthly data (~96 observations) with an additive model (period=12). However, there are a few issues and potential improvements to ensure robustness, clarity, and correctness, especially given your goal of confirming the strength of the seasonal component. Below, I'll evaluate your code, identify issues, suggest improvements, and confirm whether it will work for quantifying seasonality strength, keeping the response concise and actionable.

**Does the Code Work?**

Yes, the code will work to quantify seasonality strength by calculating the proportion of variance explained by the seasonal component (np.var(result_monthly.seasonal) / np.var(df_monthly['count'])), provided the data is correctly formatted and df_monthly is defined. However, there are critical issues and inefficiencies that could cause errors or misinterpretations, particularly around data preparation and variable usage. Let's break it down.

**Issues in the Code**

1. **Undefined df_monthly**:
   - In the final block, you use df_monthly["count"] for decomposition and variance calculation, but df_monthly is not defined anywhere in the code. You define df (from parquet) and p (grouped by year/month), but df_monthly appears to be a typo or missing variable.
   - **Impact**: This will raise a NameError and prevent the seasonality strength calculation from running.
   - **Fix**: Replace df_monthly with p (the monthly aggregated DataFrame with a count column and datetime index), as it's the series you decomposed earlier.

2. **Redundant Decomposition**:

- o You perform decomposition twice: once on p["count"] for plotting (result) and again on df_monthly["count"] for variance calculation (result_monthly). This is inefficient and risks inconsistency if df_monthly were defined differently.

  - o **Impact**: Wastes computation and could lead to confusion if results differ.

  - o **Fix**: Reuse the result from the first decomposition for variance calculation.

3. **Data Preparation Complexity**:

   - o The grouping and reindexing to ensure all months (1-12) for each year are present is correct but overly complex. Creating all_months and looping over years can be simplified using pandas' date range functionality.

   - o **Impact**: The code works but is harder to maintain and less efficient.

   - o **Fix**: Use pd.date_range to create a complete monthly index and reindex p directly.

4. **Dropping First Row**:

   - o You drop the first row of df (df.drop(index=[0])) without explanation. This could be intentional (e.g., removing an invalid row), but if it's unnecessary, it may discard valid data.

   - o **Impact**: Potential data loss if not justified.

   - o **Fix**: Clarify or remove this step unless it's specific to your data's structure.

5. **Error Handling for Decomposition**:

   - o The check if len(p) < 12 is good but may be too restrictive. Decomposition can work with fewer than 12 periods if the data is sufficient, though results may be less reliable.

   - o **Impact**: Could prematurely exit for valid datasets.

   - o **Fix**: Adjust to len(p) < 2 * period (e.g., 24 for period=12) to ensure at least two seasonal cycles, as recommended for stable decomposition.

6. **Logging and Plotting**:

   - o The logging setup with loguru is robust, but the plot saving (plt.savefig) occurs before plt.show(), which is fine but could be streamlined.

   - o **Impact**: Minor inefficiency in plot handling.

   - o **Fix**: Combine plotting and saving in a single block and ensure plt.close() to avoid memory issues.

7. **Configuration Dependency**:

   o The code relies on a config.toml file and a specific data/processed directory structure. If these are missing or misconfigured, the code fails early.

   o **Impact**: Reduces portability for others or different environments.

   o **Fix**: Add fallback options or validate the config structure explicitly.

**Improvements to the Code**

Below is a revised version of your code, focusing on the seasonality strength calculation for monthly data, addressing the issues, and incorporating best practices. I've streamlined data preparation, fixed the df_monthly issue, and optimized for clarity and efficiency.

**Key Improvements**

1. **Fixed df_monthly**:

   o Replaced df_monthly with monthly (the aggregated series), ensuring the variance calculation uses the same data as the decomposition.

2. **Simplified Data Preparation**:

   o Used pd.date_range and reindex to fill missing months efficiently, replacing the complex all_months loop.

   o Aggregated directly to monthly counts using to_period("M") for cleaner grouping.

3. **Single Decomposition**:

   o Performed decomposition once (result) and reused it for both plotting and variance calculation, improving efficiency.

4. **Adjusted Length Check**:

   o Changed len(p) < 12 to len(monthly) < 24 to ensure at least two seasonal cycles (2 × 12 = 24 months), aligning with best practices for stable decomposition.

5. **Plotting Cleanup**:

   o Added plt.close() to free memory and streamlined plot saving.

6. **Removed Unnecessary Drop**:

   o Omitted df.drop(index=[0]) since it's unclear why it's needed. If specific to your data, add it back with a comment explaining why.

7. **Robust Logging**:

   o Kept your logging structure but added more specific messages for clarity (e.g., data preparation steps).

**Will It Work?**

- **Yes**: The revised code will correctly quantify seasonality strength by calculating np.var(result.seasonal) / np.var(monthly), assuming your data is in the expected format (parquet file with a timestamp column). It produces a decomposition plot and logs the variance proportion, confirming the strength of the monthly seasonal component (e.g., >20% indicates strong seasonality).

- **Output**: You'll get a log entry like "Monthly Seasonal Variance Proportion: 0.250" and a decomposition plot saved to img/decomposition-month.png.

**Additional Suggestions**

1. **Validate Input Data**:

   o Add checks for the timestamp column and non-negative counts:

```
if "timestamp" not in df.columns:

  logger.error("DataFrame missing 'timestamp' column")

  exit(1)

if (df.groupby("year-month").size() < 0).any():

  logger.error("Negative counts detected")

  exit(1)
```

2. **Compare with Weekly Data**:

   o If you want to confirm monthly data's superiority, compute the weekly seasonal variance (period=52) and compare:

```
weekly = df.groupby(df["date"].dt.isocalendar().week).size()

result_weekly = seasonal_decompose(weekly, model="additive", period=52)

weekly_variance = np.var(result_weekly.seasonal) / np.var(weekly)

logger.info(f"Weekly Seasonal Variance Proportion: {weekly_variance:.3f}")
```

3. **Residual Analysis**:

   o Check residuals to ensure the decomposition is valid:

```
from statsmodels.graphics.tsaplots import plot_acf
```

```
plot_acf(result.resid.dropna(), lags=24, title="Monthly Residual ACF")

plt.savefig("img/residual_acf.png", dpi=300, bbox_inches="tight")

plt.close()
```

- If ACF spikes exceed ±0.2 (for ~96 observations), reconsider the period or model.

**Conclusion**

Your code is on the right track but had a critical error (df_monthly undefined) and inefficiencies (redundant decomposition, complex reindexing). The revised version fixes these, streamlines data preparation, and robustly quantifies seasonality strength. It will work to confirm your observation of stronger monthly seasonality (e.g., variance proportion >20%) and produce a clear decomposition plot. Run the revised code, and if the seasonal variance is high, proceed to SARIMA modeling (as suggested previously) for forecasting.

If you share the output (e.g., seasonal variance value) or data specifics, I can further validate or suggest next steps. Would you like me to extend the code for SARIMA forecasting or generate a sample decomposition plot?

Monthly Seasonal Variance Proportion: 0.424. What does this prove?

A **Monthly Seasonal Variance Proportion of 0.424** (42.4%) from your additive decomposition of 8 years of monthly data (period=12) is a strong indicator of **significant annual seasonality** in your time series. Below, I'll explain what this proves, its implications, and what to do next, keeping the response concise and grounded in time series analysis.

**What It Proves**

1. **Strong Seasonality**:
   - The seasonal component explains 42.4% of the total variance in your monthly data, which is substantial (typically, >20% is considered strong). This proves that **annual seasonality is a major driver** of your time series' behavior, with consistent patterns repeating every 12 months (e.g., peaks in specific months like December, troughs in others like July).
   - **Context**: Your earlier finding of stronger seasonality in monthly data (vs. weekly) and low ACF values (+0.1 to -0.1 for weekly data) suggests that monthly aggregation smoothed out noise, revealing a clear annual cycle that was masked in weekly data.

2. **Additive Model Fit**:

- o Since you used an additive model (Y = Trend + Seasonal + Residual), the high seasonal variance confirms that seasonal fluctuations are **constant in magnitude** across the 8 years, independent of the trend. This validates the choice of an additive model over a multiplicative one for your data.

3. **Reliability of Decomposition**:

    - o With 8 years (~96 monthly observations), the decomposition is robust (rule of thumb: at least 2-4 cycles, you have 8). The high variance proportion indicates the decomposition successfully isolated a meaningful seasonal component, especially since monthly data aligns well with calendar-driven cycles.

**Implications**

- **Seasonality Dominates**: Nearly half the variation in your data is due to predictable monthly patterns, making seasonality critical for modeling and forecasting.

- **Monthly Data Superior**: The strong seasonality in monthly data (vs. weaker weekly results) confirms that monthly aggregation is better for capturing annual cycles, likely due to reduced noise.

- **Actionable Patterns**: Specific months likely drive the signal (e.g., holiday peaks), which can guide business decisions or further analysis.

**Next Steps**

Given the strong seasonality (42.4% variance), leverage it for modeling and forecasting:

1. **Validate Decomposition**:

    - o **Action**: Check residuals to ensure the model captured all systematic patterns:

```
from statsmodels.graphics.tsaplots import plot_acf

from statsmodels.tsa.seasonal import seasonal_decompose

result = seasonal_decompose(df_monthly['count'], model='additive', period=12)

plot_acf(result.resid.dropna(), lags=24, title="Monthly Residual ACF")

plt.savefig("img/residual_acf.png", dpi=300, bbox_inches="tight")

plt.close()

from statsmodels.stats.diagnostic import acorr_ljungbox

lb_test = acorr_ljungbox(result.resid.dropna(), lags=[12], return_df=True)
```

```
print(f"Ljung-Box p-value at lag 12: {lb_test['lb_pvalue'].iloc[0]:.3f}")
```

- o **Why**: If ACF values are within ±0.2 (~95% CI for 96 observations) and Ljung-Box p>0.05, residuals are random, confirming the decomposition's validity. If spikes remain, try period=6 (semi-annual) or a multiplicative model.

2. **Fit a Seasonal Model**:

- o **Action**: Use SARIMA to model the series, incorporating the strong seasonality (period=12):

```
from statsmodels.tsa.statespace.sarimax import SARIMAX

model = SARIMAX(df_monthly['count'], order=(1,1,1), seasonal_order=(1,1,1,12))

fit = model.fit()

forecast = fit.forecast(steps=12)  # 1-year forecast

forecast.plot(title="1-Year Monthly Forecast")

plt.savefig("img/forecast.png", dpi=300, bbox_inches="tight")

plt.close()
```

- o **Why**: SARIMA explicitly models the 42.4% seasonal variance, ideal for additive data with strong annual cycles.

3. **Evaluate Forecast Accuracy**:

- o **Action**: Split data to test the model:

```
train = df_monthly.iloc[:-12]

test = df_monthly.iloc[-12:]

model = SARIMAX(train['count'], order=(1,1,1), seasonal_order=(1,1,1,12)).fit()

forecast = model.forecast(steps=12)

mse = np.mean((forecast - test['count'])**2)

print(f"Forecast MSE: {mse:.3f}")
```

- o **Why**: Low MSE confirms the model leverages seasonality effectively. Compare with a non-seasonal ARIMA to quantify seasonality's benefit.

4. **Identify Key Months**:

- o **Action**: Analyze the seasonal component to pinpoint months driving the 42.4% variance:

```
seasonal = result.seasonal

seasonal_df = pd.DataFrame({'month': seasonal.index.month, 'seasonal':
seasonal.values})

seasonal_by_month = seasonal_df.groupby('month').mean()

seasonal_by_month.plot(kind='bar', title="Average Seasonal Effect by Month")

plt.savefig("img/seasonal_by_month.png", dpi=300, bbox_inches="tight")

plt.close()
```

- o **Why**: Identifies months with the largest seasonal effects (e.g., December peaks) for business insights.

5. **Explore External Drivers**:

- o **Action**: If residuals show patterns or forecasts underperform, regress residuals on covariates (e.g., holidays):

```
df_monthly['month'] = df_monthly.index.month

dummies = pd.get_dummies(df_monthly['month'], prefix='month')

from statsmodels.api import OLS

model = OLS(result.resid, dummies).fit()

print(model.summary())
```

- o **Why**: Pinpoints specific months or events driving seasonality.

**Chart to Visualize Seasonal Component**

To illustrate the strong seasonality (42.4% variance), here's a sample bar chart of the average seasonal effect by month, assuming typical patterns (e.g., peaks in December). Since I don't have your exact seasonal component, this is illustrative:

Grok kan fouten maken. Controleer altijd de originele bronnen.Downloaden

**Interpretation**: This chart shows a clear annual pattern (e.g., December peak, July trough), reflecting the 42.4% variance explained by seasonality.

**Conclusion**

- **What It Proves**: The 42.4% seasonal variance proportion confirms **strong annual seasonality** in your monthly data, with consistent patterns explaining nearly half the variation. The additive model is appropriate, and monthly data is ideal for modeling due to its clarity over weekly data.

- **Next Steps**: Validate residuals, fit a SARIMA model (period=12), evaluate forecast accuracy, and analyze key months for insights. SARIMA is recommended for leveraging the strong seasonality.

- **Action**: Run the residual ACF check and SARIMA code above to confirm the model and generate forecasts.

If you share residual ACF results or specific seasonal patterns (e.g., which months peak), I can refine the SARIMA parameters or provide a tailored forecast. Would you like a sample SARIMA forecast script or a plot based on dummy data?

Trying to validate Decomposition with code below. Apparently there is an issue with df_monthly.

Your code for validating the decomposition of monthly data by checking residual autocorrelation and performing a Ljung-Box test is on the right track, but as you noted, there's an issue with df_monthly, which is undefined and will cause a NameError. This is similar to the issue in your previous code, where df_monthly was incorrectly referenced instead of the processed monthly series (monthly). Below, I'll confirm the issue, explain why it occurs, provide a fixed version of the code, and suggest improvements to ensure robust validation of your decomposition, given the strong monthly seasonality (42.4% variance proportion) you reported. The response will be concise, actionable, and tailored to your 8 years of monthly data (~96 observations) with an additive model (period=12).

**Issue with df_monthly**

- **Problem**: The code references df_monthly['count'] in the decomposition step for residual analysis (result = seasonal_decompose(df_monthly['count'], model='additive', period=12)), but df_monthly is not defined. You create a monthly series (a pandas Series with a datetime index and count values) from aggregating df, which is the correct data for decomposition.

- **Impact**: This will raise a NameError: name 'df_monthly' is not defined, halting execution.

- **Cause**: Likely a typo or carryover from earlier code, where you intended to use monthly (the aggregated monthly counts).

- **Fix**: Replace df_monthly['count'] with monthly, as it's the correctly formatted monthly series used in the earlier decomposition.

**Other Observations**

1. **Redundant Decomposition**: You decompose the series twice: once for plotting and variance calculation (result = seasonal_decompose(monthly, ...)), and again

for residual analysis (result = seasonal_decompose(df_monthly['count'], …)). This is inefficient and risks inconsistency if the data differed.

2. **Robust Structure**: The code's data preparation, logging, and error handling are solid, and the residual validation steps (ACF plot and Ljung-Box test) are appropriate for confirming the decomposition's validity.

3. **Minor Improvements**: The Ljung-Box test and ACF plot are well-implemented, but you could add checks for data quality and enhance logging for clarity.

## Revised Code

Below is a corrected version of your code, fixing the df_monthly issue, removing redundant decomposition, and adding minor improvements for robustness and clarity. The code validates the decomposition by checking residual randomness via ACF and the Ljung-Box test.

## Key Fixes and Improvements

1. **Fixed df_monthly**:

   o Replaced df_monthly['count'] with monthly in the residual analysis, reusing the result from the earlier decomposition to avoid redundancy.

2. **Streamlined Decomposition**:

   o Performed decomposition once, using its results for both variance calculation and residual validation, improving efficiency.

3. **Enhanced Logging**:

   o Added a specific log message for the validation step to track progress clearly.

4. **Data Quality Check**:

   o Kept the len(monthly) < 24 check, ensuring at least two seasonal cycles (2 × 12 = 24 months) for robust decomposition, aligning with your 8 years (~96 observations).

5. **Plot Handling**:

   o Ensured plt.close() for both decomposition and ACF plots to prevent memory issues.

## Will It Work?

• **Yes**: The revised code will validate the decomposition by:

- o   Plotting the residual ACF to check for randomness (values within ±0.2, ~95% CI for ~96 observations).

- o   Running the Ljung-Box test at lag 12 to test for residual autocorrelation (p>0.05 indicates random residuals, confirming the decomposition captured the seasonality).

- o   Outputting the seasonal variance proportion (expecting ~0.424, as you reported).

- **Expected Output**:

  - o   A decomposition plot saved to img/decomposition-month.png.

  - o   A residual ACF plot saved to img/residual_acf.png.

  - o   Log entries, e.g., "Monthly Seasonal Variance Proportion: 0.424" and "Ljung-Box p-value at lag 12: 0.xxx".

## Interpreting Results

- **Residual ACF**: If most ACF values are within ±0.2 and no significant spikes occur at lag 12 or multiples, the decomposition successfully captured the seasonality (42.4% variance). If spikes exceed ±0.2, consider a different period (e.g., 6 for semi-annual) or a multiplicative model.

- **Ljung-Box Test**: A p-value >0.05 at lag 12 indicates residuals are random, validating the additive model with period=12. If p<0.05, residual autocorrelation remains, suggesting a mis-specified period or model.

## Additional Suggestions

1. **Check Data Quality**:

   - o   Add validation for the timestamp column and non-negative counts:

```
if "timestamp" not in df.columns:

  logger.error("DataFrame missing 'timestamp' column")

  exit(1)

if (monthly < 0).any():

  logger.error("Negative counts detected")

  exit(1)
```

2. **Test Alternative Periods**:

   - o   If residuals show patterns, try period=6 or 13:

```
result_alt = seasonal_decompose(monthly, model="additive", period=6)
```

```
logger.info(f"Residual Variance (period=6): {np.var(result_alt.resid.dropna()):.3f}")
```

3. **Proceed to Modeling**:

  o  Given the strong seasonality (42.4%), fit a SARIMA model next:

```
from statsmodels.tsa.statespace.sarimax import SARIMAX
```

```
model = SARIMAX(monthly, order=(1,1,1), seasonal_order=(1,1,1,12))
```

```
fit = model.fit()
```

```
forecast = fit.forecast(steps=12)
```

```
plt.plot(forecast, title="1-Year Monthly Forecast")
```

```
plt.savefig("img/forecast.png", dpi=300, bbox_inches="tight")
```

```
plt.close()
```

**Conclusion**

- **Issue Fixed**: The df_monthly error is resolved by using monthly and reusing the decomposition result. The code now correctly validates the decomposition.

- **Validation**: The ACF plot and Ljung-Box test will confirm if the 42.4% seasonal variance is fully captured. Expect random residuals (p>0.05) given the strong seasonality.

- **Next Step**: Run the revised code. If residuals are random, proceed to SARIMA modeling for forecasting, leveraging the strong monthly seasonality.

If you share the Ljung-Box p-value or ACF plot results, I can interpret them or suggest further refinements. Would you like a sample SARIMA forecast script or a mock ACF plot for random residuals?

Ljung-Box p-value at lag 12: 0.086 What does this imply?

A **Ljung-Box p-value of 0.086** at lag 12 for the residuals of your additive decomposition of monthly data (period=12, 8 years, ~96 observations) indicates that **the residuals are likely random**, supporting the validity of your decomposition and confirming the strong seasonality (42.4% variance proportion) you identified. Below, I'll explain what this p-value implies, its significance in the context of your analysis, and the next steps, keeping the response concise and grounded in time series analysis.

**What the Ljung-Box p-value of 0.086 Implies**

1. **Residuals Are Random**:

- The Ljung-Box test assesses whether residuals exhibit significant autocorrelation at specified lags (here, lag 12, corresponding to one annual cycle). A p-value of 0.086 (>0.05) fails to reject the null hypothesis that the residuals are white noise (no autocorrelation).

- **Implication**: Your additive decomposition with period=12 successfully captured the seasonal and trend components, leaving residuals with no significant annual patterns. The 42.4% seasonal variance is reliable, and the model is well-specified for monthly data.

2. **Strong Seasonality Validated**:

- Since the residuals lack significant autocorrelation at lag 12, the strong seasonality (42.4% variance) is fully accounted for by the decomposition. This confirms that monthly data's annual cycle (period=12) is a major driver, unlike the weaker weekly data where noise obscured seasonality (low ACF values of +0.1 to -0.1).

3. **Borderline Significance**:

- A p-value of 0.086 is close to the 0.05 threshold, suggesting a small chance of minor residual autocorrelation. While not statistically significant at the 5% level, it's worth checking the residual ACF plot to ensure no subtle patterns remain.

**Contextual Implications**

- **Model Fit**: The additive model (Y = Trend + Seasonal + Residual) with period=12 is appropriate, as it captures the strong seasonality without leaving systematic patterns in the residuals.

- **Monthly Data Superiority**: The random residuals and high seasonal variance (42.4%) reinforce that monthly aggregation is better than weekly for modeling, as it reduces noise and highlights clear annual patterns.

- **Forecasting Readiness**: The validated decomposition supports using a seasonal model (e.g., SARIMA) for accurate forecasting, leveraging the strong monthly seasonality.

**Next Steps**

Given the Ljung-Box p-value of 0.086 and strong seasonality, proceed with modeling and further validation:

1. **Inspect Residual ACF Plot**:

- **Action**: Review the residual ACF plot (saved to img/residual_acf.png in your code) to confirm no significant spikes at lag 12 or other lags:

```
from statsmodels.graphics.tsaplots import plot_acf

from statsmodels.tsa.seasonal import seasonal_decompose

result = seasonal_decompose(monthly, model="additive", period=12)

plot_acf(result.resid.dropna(), lags=24, title="Monthly Residual ACF")

plt.savefig("img/residual_acf.png", dpi=300, bbox_inches="tight")

plt.close()
```

- o **Why**: Ensure ACF values are within ±0.2 (~95% CI for ~96 observations). If spikes approach or exceed this at lag 12 or multiples, consider a multiplicative model or different period (e.g., 6 for semi-annual).

2. **Fit a SARIMA Model**:

- o **Action**: Use SARIMA to model the series, leveraging the validated seasonality:

```
from statsmodels.tsa.statespace.sarimax import SARIMAX

import pandas as pd

# Assuming 'monthly' is your monthly series (as in your code)

model = SARIMAX(monthly, order=(1,1,1), seasonal_order=(1,1,1,12))

fit = model.fit()

forecast = fit.forecast(steps=12)  # 1-year forecast

plt.plot(monthly.index[-24:], monthly[-24:], label="Observed")

plt.plot(forecast.index, forecast, label="Forecast")

plt.title("1-Year Monthly Forecast")

plt.legend()

plt.savefig("img/forecast.png", dpi=300, bbox_inches="tight")

plt.close()
```

- o **Why**: SARIMA explicitly models the 42.4% seasonal variance, ideal for additive data with random residuals.

3. **Evaluate Forecast Accuracy**:

- o **Action**: Test the model on a holdout set (e.g., last 12 months):

```
train = monthly.iloc[:-12]
```

```
test = monthly.iloc[-12:]
```

```
model = SARIMAX(train, order=(1,1,1), seasonal_order=(1,1,1,12)).fit()
```

```
forecast = model.forecast(steps=12)
```

```
mse = np.mean((forecast - test)**2)
```

```
print(f"Forecast MSE: {mse:.3f}")
```

- o **Why**: Low MSE confirms the model leverages seasonality effectively. Compare with a non-seasonal ARIMA to quantify the benefit of seasonality.

4. **Analyze Seasonal Patterns**:

- o **Action**: Identify months driving the 42.4% variance:

```
seasonal = result.seasonal
```

```
seasonal_df = pd.DataFrame({'month': seasonal.index.month, 'seasonal': seasonal.values})
```

```
seasonal_by_month = seasonal_df.groupby('month').mean()
```

```
seasonal_by_month.plot(kind='bar', title="Average Seasonal Effect by Month")
```

```
plt.savefig("img/seasonal_by_month.png", dpi=300, bbox_inches="tight")
```

```
plt.close()
```

- o **Why**: Highlights key months (e.g., December peaks) for business insights.

5. **Check for Covariates**:

- o **Action**: If forecasts underperform or residuals show minor patterns, regress residuals on external factors (e.g., holidays):

```
monthly_df = pd.DataFrame({'count': monthly, 'month': monthly.index.month})
```

```
dummies = pd.get_dummies(monthly_df['month'], prefix='month')
```

```
from statsmodels.api import OLS
```

```
model = OLS(result.resid, dummies).fit()
```

```
print(model.summary())
```

- o **Why**: Identifies specific months or events driving seasonality.

**Chart to Visualize Residual ACF**

To illustrate what a random residual ACF plot should look like (given $p=0.086$), here's a sample chart for residuals with no significant autocorrelation:

**Interpretation**: This chart shows ACF values within ±0.2, with no significant spike at lag 12, consistent with your p=0.086, confirming random residuals.

**Conclusion**

- **Implication**: The Ljung-Box p-value of 0.086 indicates random residuals, validating your additive decomposition with period=12. The 42.4% seasonal variance is reliable, proving strong annual seasonality in monthly data.

- **Next Steps**: Fit a SARIMA model (period=12) for forecasting, evaluate accuracy, and analyze seasonal patterns by month. Check the residual ACF plot to ensure no minor patterns (given the borderline p-value).

- **Action**: Run the SARIMA code above and review the ACF plot. If spikes are within ±0.2, proceed with forecasting.

If you share the residual ACF plot results or specific data details, I can refine the SARIMA parameters or interpret further. Would you like a sample SARIMA forecast plot or additional code for seasonal pattern analysis?

Explain the residual_acf.png. On the x-axis is from 0 to 24 and the y-axis is from +1,00 to -1.00. All values are kind of between +0.2 and -0.2.

The residual_acf.png plot you described, generated from your code to validate the additive decomposition of your monthly data (8 years, ~96 observations, period=12), is an **Autocorrelation Function (ACF) plot of the residuals**. With a Ljung-Box p-value of 0.086 at lag 12 and a seasonal variance proportion of 42.4%, this plot helps confirm whether your decomposition successfully captured the strong annual seasonality, leaving random residuals. Below, I'll explain the plot's components, interpret the values, and outline implications and next steps, keeping the explanation concise and grounded in time series analysis.

**Explanation of the Residual ACF Plot**

1. **X-Axis (Lags, 0 to 24)**:

   - The x-axis represents the **lags** (time offsets) at which autocorrelation is calculated, ranging from 0 to 24 months. Lag 0 is the correlation of the residuals with themselves (always 1), lag 1 is the correlation with the previous month, lag 12 is one annual cycle, and so on.

   - **Why 0 to 24?**: Your code specifies lags=24 in plot_acf(result.resid.dropna(), lags=24), covering two annual cycles (2 × 12 = 24 months), which is sufficient to detect any remaining seasonal patterns in residuals for monthly data with period=12.

2. **Y-Axis (Autocorrelation, +1.00 to -1.00)**:

   o The y-axis shows the **autocorrelation coefficients** for each lag, ranging from -1 (perfect negative correlation) to +1 (perfect positive correlation). A value of 0 indicates no correlation.

   o **Interpretation**: Autocorrelation measures how similar the residuals are to their past values at each lag. Random residuals (white noise) should have ACF values near 0, within the confidence interval.

3. **Values Between +0.2 and -0.2**:

   o You noted that all ACF values (except possibly lag 0, which is always 1) are between +0.2 and -0.2. This is within or close to the **95% confidence interval** for random residuals, calculated as ±1.96/√n ≈ ±0.2 for n≈96 observations.

   o **Implication**: Values within ±0.2 suggest **no significant autocorrelation**, meaning the residuals are likely random (white noise). This confirms that the additive decomposition (period=12) captured the strong seasonality (42.4% variance) and trend, leaving no systematic patterns in the residuals.

4. **Confidence Interval**:

   o The plot likely includes shaded or dashed lines at ±0.2, representing the 95% confidence interval for a white noise process. If all bars (except lag 0) are within these bounds, it supports the conclusion that residuals lack significant autocorrelation.

**What This Proves**

- **Valid Decomposition**: The ACF values between +0.2 and -0.2, combined with the Ljung-Box p-value of 0.086 (>0.05) at lag 12, confirm that the residuals are random. This validates your additive decomposition, indicating that the 42.4% seasonal variance and trend components fully capture the systematic patterns in your monthly data.

- **Strong Seasonality Confirmed**: The absence of significant residual autocorrelation at lag 12 (or multiples like 24) proves that the annual seasonality (period=12) is well-modeled, aligning with the strong seasonal variance you reported.

- **No Missed Patterns**: The lack of spikes beyond ±0.2 suggests no remaining annual or other periodic patterns, reinforcing that monthly data's seasonality is stronger and clearer than weekly data (where low ACF values indicated noise).

**Implications**

- **Model Fit**: The additive model with period=12 is appropriate, as it leaves random residuals, supporting its use for forecasting or analysis.

- **Monthly Data Advantage**: The clear seasonality (42.4% variance) and random residuals in monthly data confirm it's superior to weekly data for modeling, as it reduces noise and highlights annual cycles.

- **Forecasting Readiness**: The validated decomposition supports using a seasonal model like SARIMA to leverage the strong seasonality for accurate predictions.

**Potential Caveats**

- **Borderline p-value**: The Ljung-Box p-value of 0.086 is close to 0.05, so minor autocorrelation could exist. If any ACF values are close to ±0.2 (e.g., at lag 12 or 24), it's worth investigating further to ensure no subtle patterns remain.

- **Data Quality**: If your data has gaps or outliers (filled with 0s in your code), these could affect residual randomness. Ensure the monthly series is clean and complete.

**Next Steps**

Given the random residuals and strong seasonality, proceed with modeling and analysis:

1. **Confirm ACF Plot Details**:

   - **Action**: Double-check the ACF plot for any values near ±0.2, especially at lags 12 or 24. If any are borderline, test an alternative period (e.g., 6 for semi-annual) or a multiplicative model:

```python
from statsmodels.tsa.seasonal import seasonal_decompose

result_alt = seasonal_decompose(monthly, model="multiplicative", period=12)

from statsmodels.graphics.tsaplots import import plot_acf

plot_acf(result_alt.resid.dropna(), lags=24, title="Multiplicative Residual ACF")

plt.savefig("img/residual_acf_multiplicative.png", dpi=300, bbox_inches="tight")

plt.close()
```

   - **Why**: Ensures no better model exists, especially since p=0.086 is borderline.

2. **Fit a SARIMA Model**:

   - **Action**: Model the series with SARIMA, leveraging the validated seasonality:

```
from statsmodels.tsa.statespace.sarimax import SARIMAX

import pandas as pd

import matplotlib.pyplot as plt

# Assuming 'monthly' is your monthly series (from your code)

model = SARIMAX(monthly, order=(1,1,1), seasonal_order=(1,1,1,12))

fit = model.fit()

forecast = fit.forecast(steps=12)  # 1-year forecast

plt.plot(monthly.index[-24:], monthly[-24:], label="Observed")

plt.plot(forecast.index, forecast, label="Forecast")

plt.title("1-Year Monthly Forecast")

plt.legend()

plt.savefig("img/forecast.png", dpi=300, bbox_inches="tight")

plt.close()
```

- o **Why**: SARIMA captures the 42.4% seasonal variance for accurate forecasting.

3. **Evaluate Forecast Accuracy**:
   - o **Action**: Test the model on a holdout set:

```
train = monthly.iloc[:-12]

test = monthly.iloc[-12:]

model = SARIMAX(train, order=(1,1,1), seasonal_order=(1,1,1,12)).fit()

forecast = model.forecast(steps=12)

mse = np.mean((forecast - test)**2)

print(f"Forecast MSE: {mse:.3f}")
```

- o **Why**: Confirms the model's predictive power using the validated seasonality.

4. **Analyze Seasonal Patterns**:
   - o **Action**: Identify months driving the 42.4% variance:

```
result = seasonal_decompose(monthly, model="additive", period=12)

seasonal = result.seasonal
```

```
seasonal_df = pd.DataFrame({'month': seasonal.index.month, 'seasonal':
seasonal.values})

seasonal_by_month = seasonal_df.groupby('month').mean()

seasonal_by_month.plot(kind='bar', title="Average Seasonal Effect by Month")

plt.savefig("img/seasonal_by_month.png", dpi=300, bbox_inches="tight")

plt.close()
```

  o **Why**: Highlights key months (e.g., December peaks) for insights.

**Sample ACF Plot Interpretation**

Your description matches a typical random residual ACF plot. Here's a chart replicating your description (values between +0.2 and -0.2, lags 0-24, y-axis -1 to +1):

Grok kan fouten maken. Controleer altijd de originele bronnen.Downloaden

**Interpretation**: The chart shows ACF values within ±0.2 for lags 1-24, with no significant spikes at lag 12 or 24, confirming random residuals and a valid decomposition.

**Conclusion**

- **What It Implies**: The residual ACF plot (values between +0.2 and -0.2) and Ljung-Box p-value of 0.086 confirm that your additive decomposition (period=12) captures the strong seasonality (42.4% variance), leaving random residuals. This validates the model and proves robust annual seasonality in monthly data.

- **Next Steps**: Fit a SARIMA model (period=12) for forecasting, evaluate accuracy, and analyze seasonal patterns by month. The borderline p-value suggests checking the ACF plot for near-threshold values, but the current model is likely sufficient.

- **Action**: Run the SARIMA code above and review the seasonal pattern analysis to identify key months.

If you can confirm specific ACF values (e.g., any close to ±0.2) or share additional data details, I can refine the next steps. Would you like a sample SARIMA forecast plot or further code for seasonal analysis?