Modeling and Forecasting Seasonality



Zhenhao Gong University of Connecticut

This course is designed to be:

- 1. Introductory
- 2. Leading by interesting questions and applications
- 3. Less math, useful, and fun!

Most important:

Feel free to ask any questions! \bigcirc

Enjoy! 🕲



▶ Unobserved Components: trend, seasonal, cycle, noise.

$$y_t = T_t + S_t + C_t + \varepsilon_t.$$

Or

$$y_t = T_t \times S_t \times C_t \times \varepsilon_t.$$

▶ We focus on seasonal on this lecture.



A seasonal pattern is one that repeats itself every year. It arises from links of technologies, preferences and institutions to the calendar.

We focus on the deterministic seasonality in which the annual repetition can be **exact**.



Example: monthly U.S. current-dollar liquor sales 1980.01 - 1992.01: very high in Nov. and Dec.



Example: Monthly U.S. current-dollar durable goods sales, 1980.01 - 1992.01: fall in Dec.



 Example: Monthly U.S. current-dollar gasoline sales, 1980.01 - 1992.01: higher in summertime



One way to deal with seasonality in a series is simply to remove it, and then to model and forecast the **seasonally adjusted series**.

Useful when interest centers explicitly on forecasting nonseasonal fluctuations.

Often inappropriate in business forecasting situations in which all the variation in a series are interested.



A key technique for modeling seasonality is regression on **seasonal dummies**.

Let ${\bf s}$ be the number of seasons in year. We can think it as the number of observations on a series in each year.

- s = 4 if we have quarterly data
- s = 12 if we have monthly data
- s = 52 if we have weekly data

Regression analysis can also be used when the regressor is **binary**, that is, when it takes on only two values, 0 or 1:

•
$$X = 1$$
 if small class size, $= 0$ if not

•
$$X = 1$$
 if female, $= 0$ if male

•
$$X = 1$$
 if treated (experimental drug), = 0 if not

Binary regressors are sometimes called "dummy" variables.



Interpreting regressions with a binary regressor

$$Y_i = \beta_0 + \beta_1 X_i + u_i, \quad i = 1, 2, \cdots, n$$

where X is binary $(X_i = 0 \text{ or } 1)$:

- When $X_i = 0$, $Y_i = \beta_0 + u_i$, the conditional mean of Y_i given $X_i = 0$ is β_0 . That is, $E(Y_i|X_i = 0) = \beta_0$.
- When $X_i = 1$, $Y_i = \beta_0 + \beta_1 + u_i$, the conditional mean of Y_i given $X_i = 1$ is $\beta_0 + \beta_1$. That is, $E(Y_i|X_i = 1) = \beta_0 + \beta_1$.

So:

$$\beta_1 = E(Y_i | X_i = 1) - E(Y_i | X_i = 0)$$

= population difference in group means.



Construct seasonal dummy variables:

▶ Use four seasons as an example,

$$D_1 = (1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, \cdots)$$

$$D_2 = (0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, \cdots)$$

$$D_3 = (0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, \cdots)$$

$$D_4 = (0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, \cdots).$$

D₁ indicates whether we're in the first quarter and so on.
At any given time, we can be in only one of the four quarters, so one seasonal dummy is 1, and all others are zero.



The pure seasonal dummy model is

$$\mathbf{y}_{\mathrm{t}} = \sum_{\mathrm{i}=1}^{\mathrm{s}} \gamma_{\mathrm{i}} \mathbf{D}_{\mathrm{it}} + \boldsymbol{\varepsilon}_{\mathrm{t}}.$$

- ▶ We're just regressing on an intercept, but allow for a different intercept in each season.
- ► Those different intercepts are called the **seasonal factors**; they summarize the seasonal pattern over the year.
- ▶ In the absence of seasonality, we can drop all the seasonal dummies and instead simply include an intercept in the usual way.



Instead of including a full set of **s** seasonal dummies, we can include any (s-1) seasonal dummies and an intercept.

- ▶ the constant term is the intercept for the omitted season.
- ▶ the coefficients on the seasonal dummies give the seasonal increase or decrease relative to the omitted season.

In no case, however, should we include s seasonal dummies and an intercept.

 cause the issue of perfect multicollinearity (dummy variable trap) if we do so.



Trend may be included as well, in which case the model is

$$y_t = \beta_1 TIME_t + \sum_{i=1}^{s} \gamma_i D_{it} + \varepsilon_t.$$

- Generalize what we did in modeling trend.
- ▶ The idea of seasonality may be extended to allow for more general calendar effects, such as holiday variation and trading-day variation.



Holiday variation refers to the fact that some holidays' dates change over time.

- ► Arrive at approximately the same time each year, the exact dates differ. (Easter)
- ▶ The behavior of many series depends in part on the timing of such holidays.
- ▶ As with seasonality, holiday effects may be handled with dummy variables.

Example: In a monthly model, in addition to a full set of seasonal dummies, we might include an "Easter dummy," which is 1 if the month contains Easter and 0 otherwise.



Trading-day variation refers to the fact that different months contain different numbers of trading days or business days.

Example: In a monthly forecasting model of volume traded on the London Stock Exchange, in addition to a full set of seasonal dummies, we might include a trading day variable, whose value each month is the number of trading days that month.



Allowing for the possibility of holiday or trading day variation gives the complete model

$$y_t = \beta_1 \text{TIME}_t + \sum_{i=1}^s \gamma_i D_{it} + \sum_{i=1}^{v_1} \delta_i^{\text{HD}} \text{HDV}_{it} + \sum_{i=1}^{v_2} \delta_i^{\text{TD}} \text{TDV}_{it} + \varepsilon_t.$$

- ▶ The HDVs are the relevant holiday variables (there are v_1 of them).
- The TDVs are the relevant trading day variables (there are v_2 of them).
- ▶ This is just a standard regression equation and can be estimated by ordinary least squares.



We consider constructing an h-step-ahead point forecast, $y_{T+h,T}$, at time T:

$$\begin{split} \mathbf{y}_{\mathrm{T+h,T}} &= \beta_{1}\mathrm{TIME_{T+h}} + \sum_{i=1}^{s}\gamma_{i}\mathbf{D}_{i,\mathrm{T+h}} + \sum_{i=1}^{v_{1}}\delta_{i}^{\mathrm{HD}}\mathrm{HDV_{it}} \\ &+ \sum_{i=1}^{v_{2}}\delta_{i}^{\mathrm{TD}}\mathrm{TDV_{it}}. \end{split}$$

Interval forecast: $y_{T+h,T} \pm 1.96\sigma$, assuming $\varepsilon_t \sim N(0, \sigma^2)$. Density forecast: $N(y_{T+h,T}, \sigma^2)$.



To make the point forecast operational, we can replace the unknown parameters with estimates:

$$\begin{split} \hat{y}_{T+h,T} &= \hat{\beta}_1 TIME_{T+h} + \sum_{i=1}^{s} \hat{\gamma}_i D_{i,T+h} + \sum_{i=1}^{v_1} \hat{\delta}_i^{HD} HDV_{i,T+h} \\ &+ \sum_{i=1}^{v_2} \hat{\delta}_i^{TD} TDV_{i,T+h}. \end{split}$$

Interval forecast: $\hat{y}_{T+h,T} \pm 1.96\hat{\sigma}$. Density forecast: $N(\hat{y}_{T+h,T}, \hat{\sigma}^2)$.



Housing starts are **seasonal** because it's usually preferable to start houses in the spring, so that they're completed before winter arrives.

Date:

Monthly data on U.S. housing starts 1946.01 - 1994.11; Use the 1946.01 - 1993.12 period for estimation; 1994.01 - 1994.11 period for out-of-sample forecasting.



▶ Housing Starts, 1946.01 - 1994.11





▶ Housing Starts, 1990.01 - 1994.11



Time



The figures reveal that there is no trend, so we'll work with the pure seasonal model,

$$\mathbf{y}_{\mathbf{t}} = \sum_{\mathbf{i}=1}^{\mathbf{s}} \gamma_{\mathbf{i}} \mathbf{D}_{\mathbf{i}\mathbf{t}} + \boldsymbol{\varepsilon}_{\mathbf{t}}.$$

Here we let s = 12, which means we have twelve seasonal dummies in the forecasting model.



LS // Dependent Variable is STARTS Sample: 1946:01 1993:12 Included observations: 576

Coefficient	Std. Error	t-Statistic	Prob.
86 50417	4 029055	21 47009	0.0000
89.50417	4.029055	22.21468	0.0000
122.8833	4.029055	30.49929	0.0000
142.1687	4.029055	35.28588	0.0000
147.5000	4.029055	36.60908	0.0000
145.9979	4.029055	36.23627	0.0000
139.1125	4.029055	34.52733	0.0000
138.4167	4.029055	34.35462	0.0000
130.5625	4.029055	32.40524	0.0000
134.0917	4.029055	33.28117	0.0000
111.8333	4.029055	27.75671	0.0000
92.15833	4.029055	22.87344	0.0000
	Coefficient 86.50417 89.50417 122.8833 142.1687 147.5000 145.9979 139.1125 138.4167 130.5625 134.0917 111.8333 92.15833	Coefficient Std. Error 86.50417 4.029055 89.50417 4.029055 122.8833 4.029055 142.1687 4.029055 147.5000 4.029055 145.9979 4.029055 139.1125 4.029055 130.5625 4.029055 134.0917 4.029055 134.0917 4.029055 111.8333 4.029055 92.15833 4.029055	CoefficientStd. Errort-Statistic86.504174.02905521.4700989.504174.02905522.21468122.88334.02905530.49929142.16874.02905535.28588147.50004.02905536.60908145.99794.02905536.23627139.11254.02905534.35462130.56254.02905532.40524134.09174.02905533.28117111.83334.02905527.7567192.158334.02905522.87344

R-squared	0.383780
Adjusted R-squared	0.371762
S.E. of regression	27.91411
Sum squared resid	439467.5
Log likelihood	-2728.825
Durbin-Watson stat	0.154140

Mean dependent var	123.3944
S.D. dependent var	35.21775
Akaike info criterion	6.678878
Schwarz criterion	6.769630
F-statistic	31.93250
Prob(F-statistic)	0.000000





- ▶ The twelve seasonal dummies account for more than a third of the variation in housing starts, as adjusted $R^2 = .371$.
- At least some of the remaining variation is cyclical, which the model is not designed to capture. (Very low Durbin-Watson statistic.)
- Rigid seasonal pattern (there's nothing in the model other than deterministic seasonal dummies) picks up a lot of the variation in housing starts.



 Estimated Seasonal Factors for Housing Starts (Twelve estimated coefficients)



Housing Starts History, 1990.01-1993.12; Forecast, 1994.01-1994.11



► Housing Starts

History, 1990.01-1993.12; Forecast, 1994.01-1994.11; Realization, 1994.01-1994.11;

