



ACEMS Forecasting Workshop

Rob J Hyndman

1 Forecast Evaluation

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

robjhyndman.com/acemsforecasting2018

- Slides
- Exercises
- Textbook
- Useful links

Key reference

Hyndman, R. J. & Athanasopoulos, G.
(2018) *Forecasting: principles and practice*, 2nd ed.

OTexts.org/fpp2/

- Free and online
- Data sets in associated R package
- R code for examples

```
install.packages("fpp2", dependencies=TRUE)
```

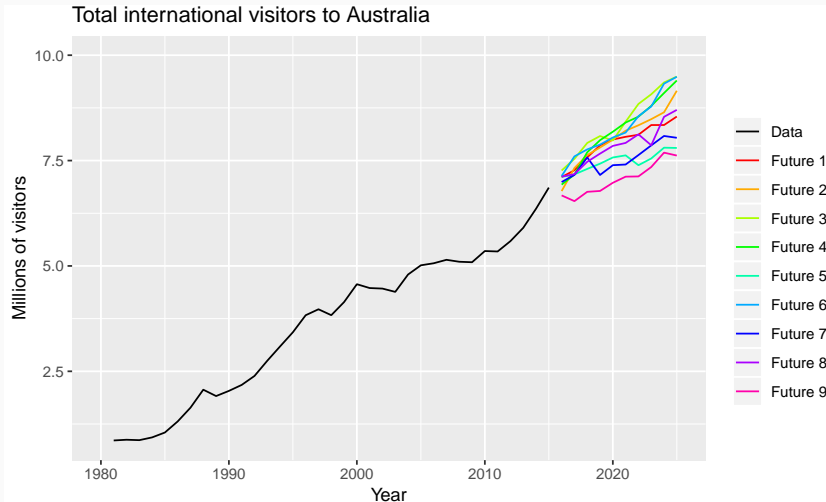
Outline

Topic	Chapter
1 Forecast evaluation	3
2 ARIMA models	8
3 Dynamic regression	9
4 Hierarchical forecasting	10
5 Ensembles	12

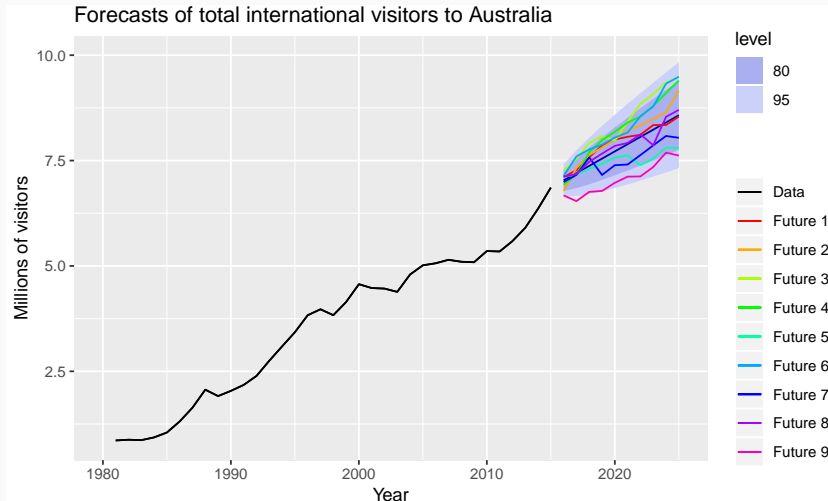
Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

Sample futures



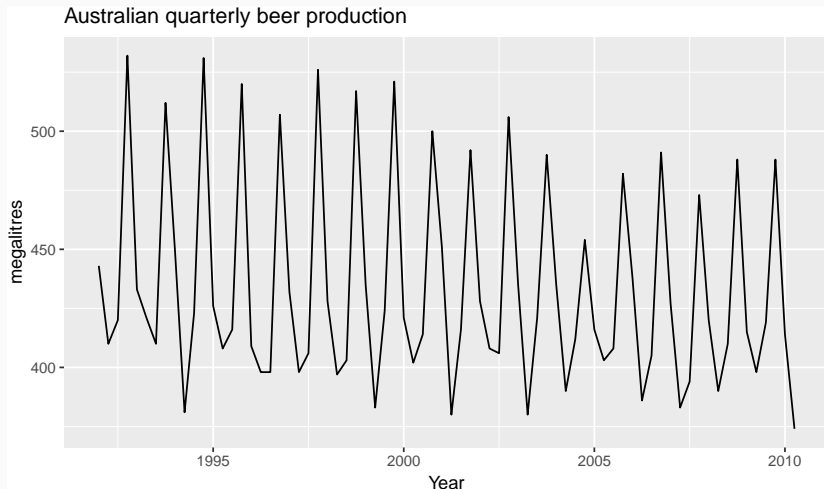
Forecast intervals



Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods**
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

Some simple forecasting methods



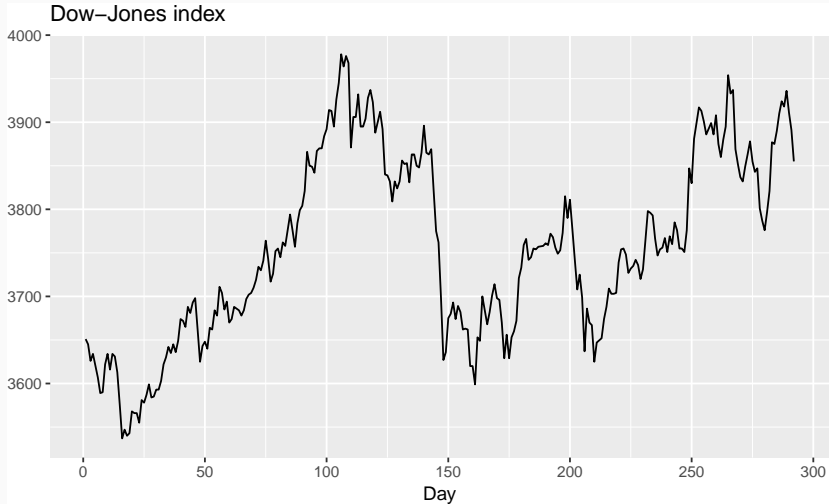
How would you forecast these data?

Some simple forecasting methods



How would you forecast these data?

Some simple forecasting methods



How would you forecast these data?

Some simple forecasting methods

Average method

- Forecast of all future values is equal to mean of historical data $\{y_1, \dots, y_T\}$.
- Forecasts: $\hat{y}_{T+h|T} = \bar{y} = (y_1 + \dots + y_T)/T$

Some simple forecasting methods

Average method

- Forecast of all future values is equal to mean of historical data $\{y_1, \dots, y_T\}$.
- Forecasts: $\hat{y}_{T+h|T} = \bar{y} = (y_1 + \dots + y_T)/T$

Naïve method

- Forecasts equal to last observed value.
- Forecasts: $\hat{y}_{T+h|T} = y_T$.
- Consequence of efficient market hypothesis.

Some simple forecasting methods

Average method

- Forecast of all future values is equal to mean of historical data $\{y_1, \dots, y_T\}$.
- Forecasts: $\hat{y}_{T+h|T} = \bar{y} = (y_1 + \dots + y_T)/T$

Naïve method

- Forecasts equal to last observed value.
- Forecasts: $\hat{y}_{T+h|T} = y_T$.
- Consequence of efficient market hypothesis.

Seasonal naïve method

- Forecasts equal to last value from same season.
- Forecasts: $\hat{y}_{T+h|T} = y_{T+h-km}$ where m = seasonal period and k is integer part of $(h - 1)/m$.

Some simple forecasting methods

Drift method

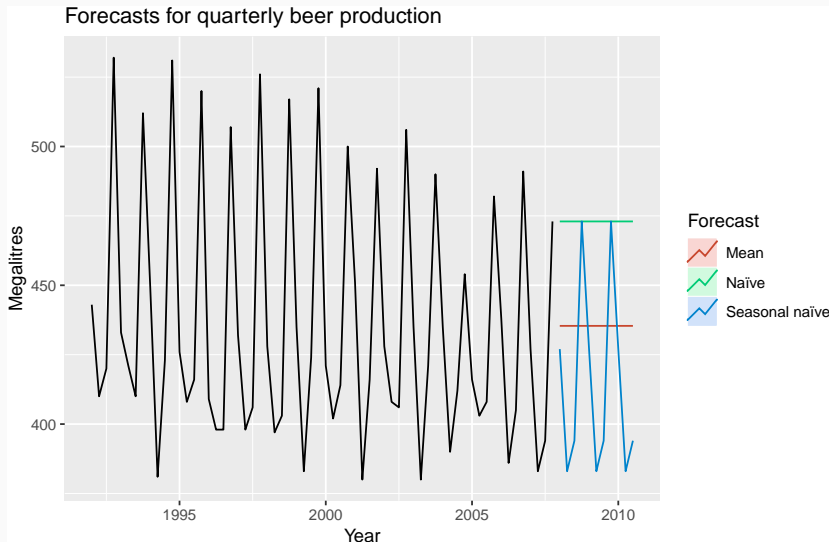
- Forecasts equal to last value plus average change.

- Forecasts:

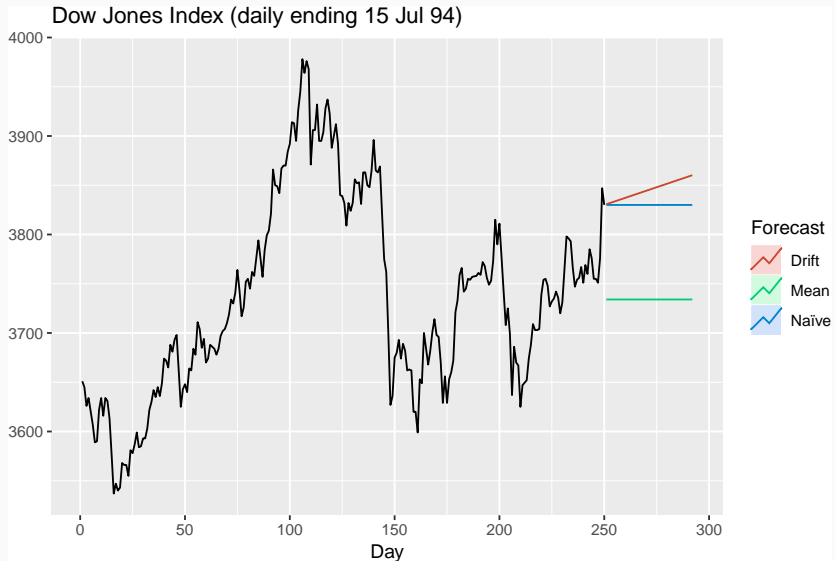
$$\begin{aligned}\hat{y}_{T+h|T} &= y_T + \frac{h}{T-1} \sum_{t=2}^T (y_t - y_{t-1}) \\ &= y_T + \frac{h}{T-1} (y_T - y_1).\end{aligned}$$

- Equivalent to extrapolating a line drawn between first and last observations.

Some simple forecasting methods



Some simple forecasting methods



Some simple forecasting methods

- Mean: `meanf(y, h=20)`
- Naïve: `naive(y, h=20)`
- Seasonal naïve: `snaive(y, h=20)`
- Drift: `rwf(y, drift=TRUE, h=20)`

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics**
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

Fitted values

- $\hat{y}_{t|t-1}$ is the forecast of y_t based on observations y_1, \dots, y_t .
- We call these “fitted values”.
- Sometimes drop the subscript: $\hat{y}_t \equiv \hat{y}_{t|t-1}$.
- Often not true forecasts since parameters are estimated on all data.

For example:

- $\hat{y}_t = \bar{y}$ for average method.
- $\hat{y}_t = y_{t-1} + (y_T - y_1)/(T - 1)$ for drift method.

Forecasting residuals

Residuals in forecasting: difference between observed value and its fitted value: $e_t = y_t - \hat{y}_{t|t-1}$.

Forecasting residuals

Residuals in forecasting: difference between observed value and its fitted value: $e_t = y_t - \hat{y}_{t|t-1}$.

Assumptions

- 1 $\{e_t\}$ uncorrelated. If they aren't, then information left in residuals that should be used in computing forecasts.
- 2 $\{e_t\}$ have mean zero. If they don't, then forecasts are biased.

Forecasting residuals

Residuals in forecasting: difference between observed value and its fitted value: $e_t = y_t - \hat{y}_{t|t-1}$.

Assumptions

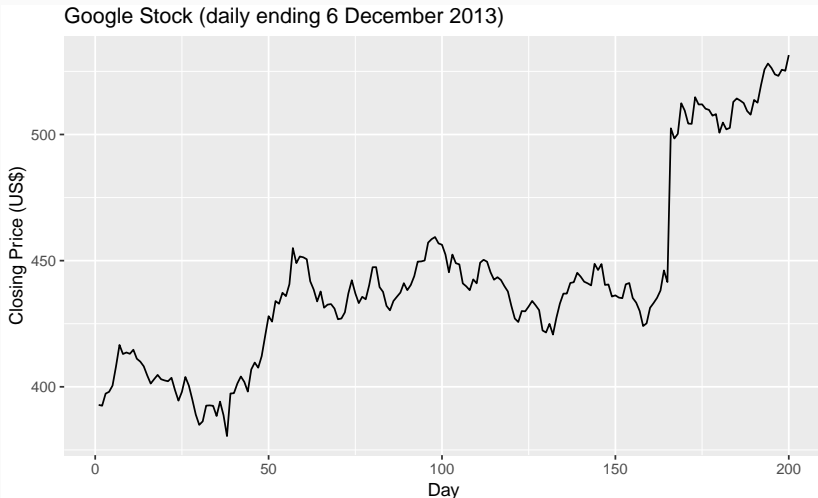
- 1 $\{e_t\}$ uncorrelated. If they aren't, then information left in residuals that should be used in computing forecasts.
- 2 $\{e_t\}$ have mean zero. If they don't, then forecasts are biased.

Useful properties (for prediction intervals)

- 3 $\{e_t\}$ have constant variance.
- 4 $\{e_t\}$ are normally distributed.

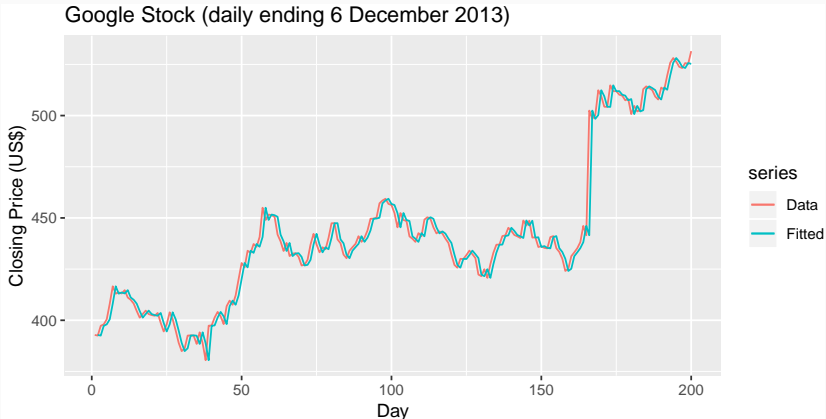
Example: Google stock price

```
autoplot(goog200) +  
  xlab("Day") + ylab("Closing Price (US$)") +  
  ggtitle("Google Stock (daily ending 6 December 2013)")
```



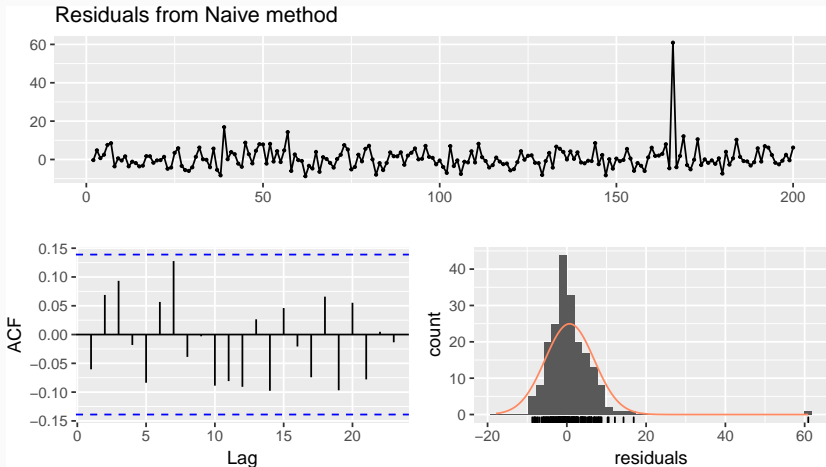
Example: Google stock price

```
fits <- fitted(naive(goog200))  
autoplot(goog200, series="Data") +  
  autolayer(fits, series="Fitted") +  
  xlab("Day") + ylab("Closing Price (US$)") +  
  ggtitle("Google Stock (daily ending 6 December 2013)")
```



checkresiduals function

```
checkresiduals(naive(goog200), test=FALSE)
```



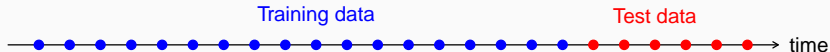
ACF of residuals

- We assume that the residuals are white noise (uncorrelated, mean zero, constant variance). If they aren't, then there is information left in the residuals that should be used in computing forecasts.
- So a standard residual diagnostic is to check the ACF of the residuals of a forecasting method.
- We *expect* these to look like white noise.

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy**
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

Training and test sets



- A model which fits the training data well will not necessarily forecast well.
- A perfect fit can always be obtained by using a model with enough parameters.
- Over-fitting a model to data is just as bad as failing to identify a systematic pattern in the data.
- The test set must not be used for *any* aspect of model development or calculation of forecasts.
- Forecast accuracy is based only on the test set.

Forecast errors

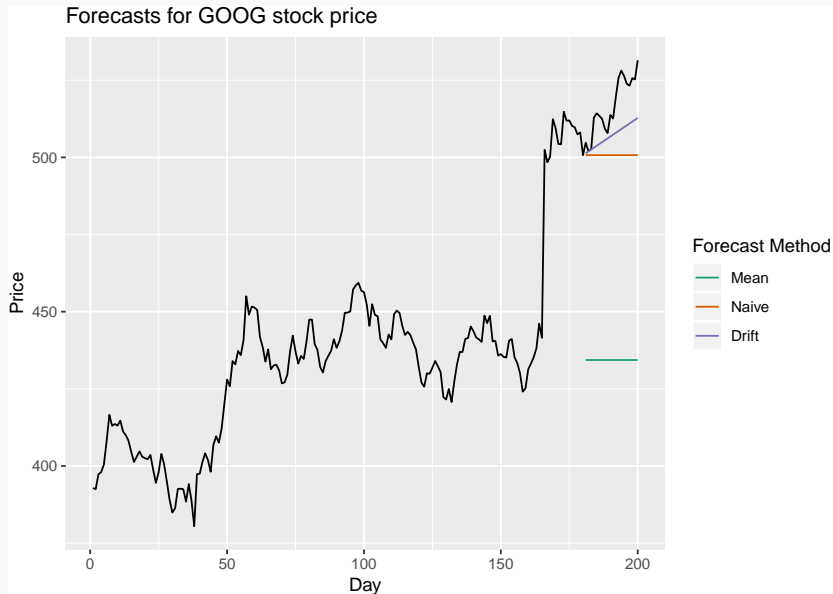
Forecast “error”: the difference between an observed value and its forecast.

$$e_{T+h} = y_{T+h} - \hat{y}_{T+h|T},$$

where the training data is given by $\{y_1, \dots, y_T\}$

- Unlike residuals, forecast errors on the test set involve multi-step forecasts.
- These are *true* forecast errors as the test data is not used in computing $\hat{y}_{T+h|T}$.

Measures of forecast accuracy



Measures of forecast accuracy

y_{T+h} = $(T + h)$ th observation, $h = 1, \dots, H$

$\hat{y}_{T+h|T}$ = its forecast based on data up to time T .

$$e_{T+h} = y_{T+h} - \hat{y}_{T+h|T}$$

$$\text{MAE} = \text{mean}(|e_{T+h}|)$$

$$\text{MSE} = \text{mean}(e_{T+h}^2)$$

$$\text{RMSE} = \sqrt{\text{mean}(e_{T+h}^2)}$$

$$\text{MAPE} = 100\text{mean}(|e_{T+h}|/|y_{T+h}|)$$

Measures of forecast accuracy

y_{T+h} = $(T + h)$ th observation, $h = 1, \dots, H$

$\hat{y}_{T+h|T}$ = its forecast based on data up to time T .

$$e_{T+h} = y_{T+h} - \hat{y}_{T+h|T}$$

$$\text{MAE} = \text{mean}(|e_{T+h}|)$$

$$\text{MSE} = \text{mean}(e_{T+h}^2)$$

$$\text{RMSE} = \sqrt{\text{mean}(e_{T+h}^2)}$$

$$\text{MAPE} = 100\text{mean}(|e_{T+h}|/|y_{T+h}|)$$

- MAE, MSE, RMSE are all scale dependent.
- MAPE is scale independent but is only sensible if $y_t \gg 0$ for all t , and y has a natural zero.

Measures of forecast accuracy

Mean Absolute Scaled Error

$$\text{MASE} = T^{-1} \sum_{t=1}^T |y_t - \hat{y}_{t|t-1}| / Q$$

where Q is a stable measure of the scale of the time series $\{y_t\}$.

Proposed by Hyndman and Koehler (IJF, 2006).

For non-seasonal time series,

$$Q = (T - 1)^{-1} \sum_{t=2}^T |y_t - y_{t-1}|$$

works well. Then MASE is equivalent to MAE relative to a naïve method.

Measures of forecast accuracy

Mean Absolute Scaled Error

$$\text{MASE} = T^{-1} \sum_{t=1}^T |y_t - \hat{y}_{t|t-1}| / Q$$

where Q is a stable measure of the scale of the time series $\{y_t\}$.

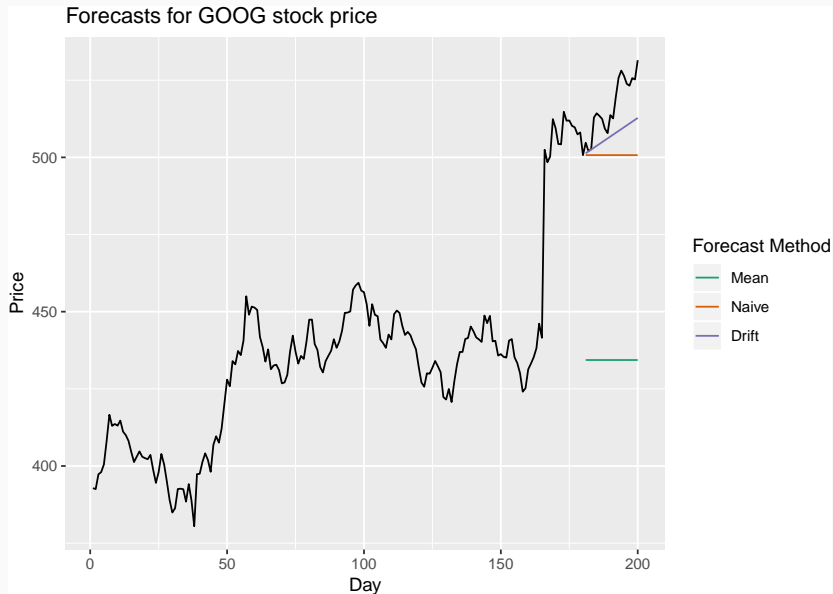
Proposed by Hyndman and Koehler (IJF, 2006).

For seasonal time series,

$$Q = (T - m)^{-1} \sum_{t=m+1}^T |y_t - y_{t-m}|$$

works well. Then MASE is equivalent to MAE relative to a seasonal naïve method.

Measures of forecast accuracy



Measures of forecast accuracy

```
googtrain <- window(goog200,end=180)
googfc1 <- meanf(googtrain,h=20)
googfc2 <- rwf(googtrain,h=20)
googfc3 <- rwf(googtrain,h=20,drift=TRUE)
accuracy(googfc1, goog200)
accuracy(googfc2, goog200)
accuracy(googfc3, goog200)
```

	RMSE	MAE	MAPE	MASE
Mean method	82.89	82.43	15.93	21.61
Naïve method	18.29	16.04	3.08	4.21
Drift method	11.34	9.71	1.86	2.55

Poll: true or false?

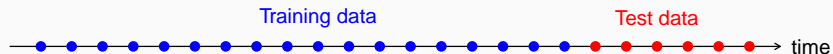
- 1 Good forecast methods should have normally distributed residuals.
- 2 A model with small residuals will give good forecasts.
- 3 The best measure of forecast accuracy is MAPE.
- 4 If your model doesn't forecast well, you should make it more complicated.
- 5 Always choose the model with the best forecast accuracy as measured on the test set.

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation**
- 7 Prediction intervals
- 8 Evaluating forecast distributions

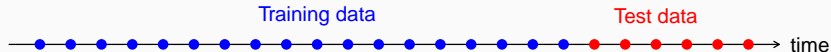
Time series cross-validation

Traditional evaluation

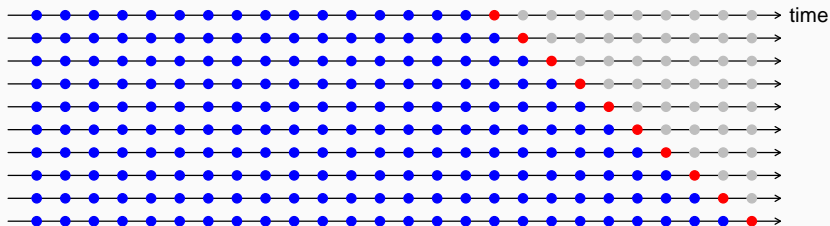


Time series cross-validation

Traditional evaluation



Time series cross-validation

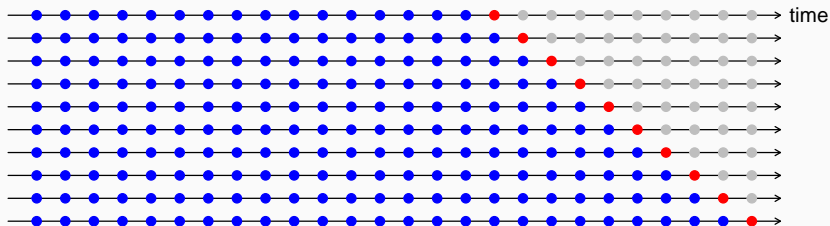


Time series cross-validation

Traditional evaluation



Time series cross-validation



- Forecast accuracy averaged over test sets.
- Also known as “evaluation on a rolling forecasting origin”

tsCV function:

```
goog200 %>%  
  tsCV(forecastfunction=meanf, h=1) -> e  
e^2 %>% mean(na.rm=TRUE) %>% sqrt
```

```
## [1] 37.16
```

```
goog200 %>%  
  tsCV(forecastfunction=naive, h=1) -> e  
e^2 %>% mean(na.rm=TRUE) %>% sqrt
```

```
## [1] 6.224
```

```
goog200 %>%  
  tsCV(forecastfunction=rwf, drift=TRUE, h=1) -> e  
e^2 %>% mean(na.rm=TRUE) %>% sqrt
```

```
## [1] 6.233
```

A good way to choose the best forecasting model is to find the model with the smallest RMSE computed using time series cross-validation.

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions

Prediction intervals

- A forecast $\hat{y}_{T+h|T}$ is (usually) the mean of the conditional distribution $y_{T+h} \mid y_1, \dots, y_T$.
- A prediction interval gives a region within which we expect y_{T+h} to lie with a specified probability.
- Assuming forecast errors are normally distributed, then a 95% PI is

$$\hat{y}_{T+h|T} \pm 1.96\hat{\sigma}_h$$

where $\hat{\sigma}_h$ is the st dev of the h -step distribution.

- When $h = 1$, $\hat{\sigma}_h$ can be estimated from the residuals.

Prediction intervals

Drift forecasts with prediction interval:

```
rwf(goog200, level=95, drift=TRUE)
```

##	Point Forecast	Lo 95	Hi 95
## 201	532.2	520.0	544.3
## 202	532.9	515.6	550.1
## 203	533.6	512.4	554.7
## 204	534.3	509.8	558.7
## 205	535.0	507.5	562.4
## 206	535.7	505.5	565.8
## 207	536.4	503.7	569.0
## 208	537.1	502.1	572.0

Prediction intervals

- Point forecasts are often useless without prediction intervals.
- Prediction intervals require a stochastic model (with random errors, etc).
- Multi-step forecasts for time series require a more sophisticated approach (with PI getting wider as the forecast horizon increases).
- Check residual assumptions before believing them.
- Usually too narrow due to unaccounted uncertainty.

Prediction intervals

Assume residuals are normal, uncorrelated, $\text{sd} = \hat{\sigma}$:

Mean forecasts: $\hat{\sigma}_h = \hat{\sigma} \sqrt{1 + 1/T}$

Naïve forecasts: $\hat{\sigma}_h = \hat{\sigma} \sqrt{h}$

Seasonal naïve forecasts $\hat{\sigma}_h = \hat{\sigma} \sqrt{k + 1}$

Drift forecasts: $\hat{\sigma}_h = \hat{\sigma} \sqrt{h(1 + h/T)}$.

where k is the integer part of $(h - 1)/m$.

Note that when $h = 1$ and T is large, these all give the same approximate value $\hat{\sigma}$.

Outline

- 1 Introduction
- 2 The statistical forecasting perspective
- 3 Benchmark methods
- 4 Residual diagnostics
- 5 Evaluating point forecast accuracy
- 6 Time series cross-validation
- 7 Prediction intervals
- 8 Evaluating forecast distributions**

Evaluating prediction intervals

Winkler score

If the $100(1 - \alpha)\%$ prediction interval is given by $[\ell, u]$, and the observed value is y , then the Winkler interval score is

$$(u - \ell) + \frac{2}{\alpha}(\ell - y)\mathbf{1}(y < \ell) + \frac{2}{\alpha}(y - u)\mathbf{1}(y > u).$$

- penalizes for wide intervals (since $u - \ell$ will be large);
- penalizes for non-coverage with observations well outside the interval being penalized more heavily.

Evaluating quantile forecasts

Let q_p be the quantile forecast with probability $1 - p$ of exceedance.

Pin-ball loss function

$$L(q_p, y) = (1 - p)(q_p - y)1(y < q_p) + p(y - q_p)1(y \geq q_p).$$

- average over all target quantiles (e.g., 0.01, 0.02, ..., 0.99) and all forecast horizons.
- Reference: Gneiting and Raftery (JASA, 2007)

Evaluating quantile forecasts