

The Art of Forecasting

Gianluca Campanella

7th June 2018

Hello!

My name is **Gianluca** [dʒan'lu:ka]

What I do nowadays

I'm a Data Scientist at



Microsoft

in Algorithms and Data Science

What I do nowadays

I also run my own company



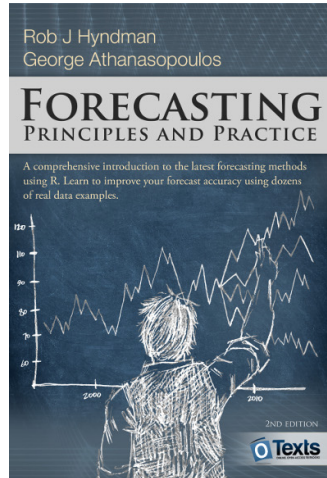
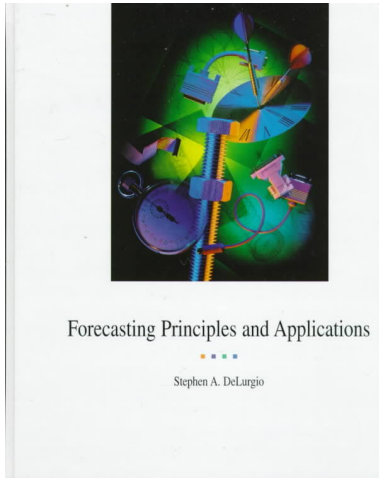
Estimand.com

that provides

Data Science training and mentoring

[https://github.com/gcampanella/
ndr-2018](https://github.com/gcampanella/ndr-2018)

References



Contents

Motivation

Modelling

Results

Recommendations

What's a time series?

Any data that change **over time**

- Typically continuous (including counts)
- Time gives natural ordering

What's forecasting?

Regression

- Value of y given values for the predictors X
- Does not depend on time (or temporal effect is negligible)

What's forecasting?

Regression

- Value of y given values for the predictors X
- Does not depend on time (or temporal effect is negligible)

Forecasting

- Value of y given **previous values** of y
- Some models can also incorporate exogenous predictors

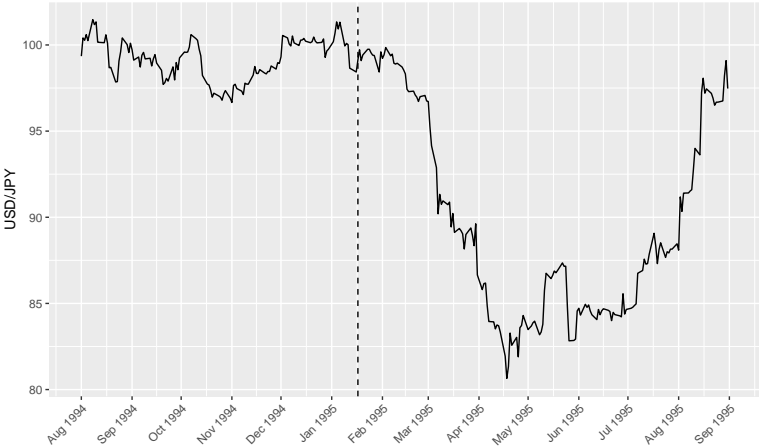
Can we forecast in changing environments?

Predictability

Predictability depends on...

- Availability of data
- Our understanding of contributing factors
- Whether our forecasts affect the process we're trying to forecast

A word of caution

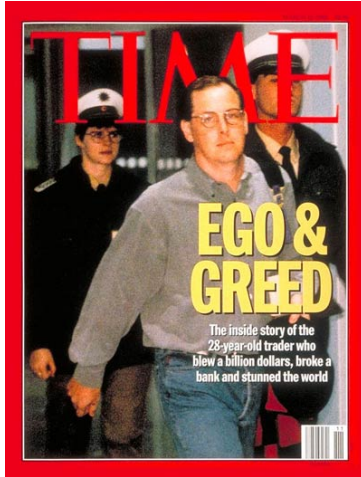


A word of caution



What happened?

A word of caution



Motivation

AVVENIRE CEI NEWS SIR TV2000 RADIO INLU FISC

Questo sito usa cookie di terze parti (anche di profilazione) e cookie tecnici. Continuando a navigare accetti l'uso. [Cambia cookie](#)

Accetta

Segui su   

 Avvenire.it

SEZIONI RUBRICHE CEI PAPA OPINIONI SINDO GIOVANI

Home - Opinioni Editoriali - Il direttore risponde - Le nostre voci

Dati come di guerra nell'Italia 2015. Attenti ai morti

Gian Carlo Bianchi venerdì 11 dicembre 2015

Leggendo i dati forniti dall'Istat sul totale dei morti in Italia nei primi sette mesi del 2015 - ultimo aggiornamento a tutt'oggi disponibile - si scopre un aumento di 39mila decessi rispetto agli stessi primi sette mesi del 2014. La cosa non è affatto marginale se si pensa che ciò corrisponde a un aumento dell'11% e che, se confermato su base annua, porterebbe a 664mila morti nel 2015 contro i 598mila dello scorso anno. Si tratterebbe di un aumento di ben 66mila unità, che si annuncia in gran parte concentrato sulla componente femminile (+40mila) e che verosimilmente coinvolgerà soprattutto la componente più anziana della popolazione residente nel nostro Paese. Il dato è impressionante. Ma ciò che lo rende del tutto anomalo è il fatto che per trovare un'analoga impennata della mortalità, con ordini di grandezza comparabili, si deve tornare indietro sino al 1943 e, prima ancora, occorre risalire agli anni tra il 1915 e il 1918: due periodi bellici della nostra storia che largamente spiegano dinamiche di questo tipo. Viceversa, in un'epoca come quella attuale, in condizioni di pace e con uno stato di benessere che, nonostante tutto, è da ritenersi ancora ampio e generalizzato, come si giustifica un rialzo della mortalità di queste dimensioni? È solo la naturale conseguenza del cambiamento in un popolo che diventa sempre più anziano o è (anche) un segnale di allarme rispetto a un sistema socio-sanitario che, dopo averci abituati al continuo allungamento della vita, - con guadagni sensibili anche in

pubblicato

OPINIONI

Secondo nel Dettling il «re delle gaffe» che alimenta anche l'euroscetticismo

Se il principe Filippo è il noto gaffeur della Casa reale britannica, il tedesco Guenther Dettling è il suo: il titolare di Blando e Biscia amaro nell'azienda.

Looking up token.rubiconproject.com...

Quando sito utilizza cookie, anche di terze parti, per inviarti pubblicità e servizi in linea con le tue preferenze. Se vuoi saperne di più o negare il consenso a tutti o ad alcuni cookie, [clicca qui](#). Chiudendo questo banner, accetti il nostro utilizzo di cookies. [OK](#)

ATTUALITÀ PARLAMENTO POLITICA POLITICA ECONOMICA DOSSIER BLOG

« [Governo, le crisi più lunghe con l'imperatore dello spread si regge l'ombra di Governo politico](#) [Sale lo spread: ecco quali sono gli effetti per le imprese e l'economia reale](#) [Spread, perché sale e perché ci deve interessare](#) [Belgio, il fallito governo](#) »

DIETRO I DATI ISTAT

In Italia nel 2015 sono morte 54mila persone in più (+9%). Ecco le possibili cause

—di Enrico Marro 25 febbraio 2016



Il rapporto Istat sugli indicatori demografici 2015 ha confermato le stime dei mesi scorsi: in Italia l'anno scorso i decessi hanno toccato quota 653mila, 54mila in più del 2014 (+9,1%). Con un tasso di mortalità, pari al 10,7 per mille, che è risultato il più alto dal secondo dopoguerra in poi. L'aumento di mortalità è concentrato tra gli anziani (75-95 anni). Come è stato possibile?

Invecchiamento popolazione e "posticipo dei decessi"
L'Istat spiega come, dal punto di vista demografico, il picco di mortalità del 2015 sia in parte dovuto a effetti strutturali connessi all'invecchiamento e in parte al posticipo delle morti non avvenute nel biennio 2013-2014. «Il picco di mortalità del 2015 porta con sé

VIDEO



30 maggio 2016
Quali effetti per i risparmiatori italiani e per le banche di crisi in caso di insuccesso della riforma? Lo spiega Maria Longo

I PIÙ LETTI DI ITALIA

- 1. COTTARELLI AL COLLE PEO «INCONTRO INFORMALE»** 29 maggio 2016
Riunione spedita di un governo M5S-Lega, Selvino e Giuglietti in pole per Palazzo Chigi. Motori aspre
- 2. CRISI ISTITUZIONALE** 30 maggio 2016
Cottarelli «C'è spazio di Governo politico, resta in attesa». Selvino: no, ma non a luglio»
- 3. DIETRO LA SVOLTA DI DI NARDI** 30 maggio 2016
M5S di lotta «dimo» di governo pena la paura di perdere il voto moderato
- 4. NONENE** 25 settembre 2015
Bernardo Mattarella nuovo ad della Banca del Mezzogiorno
- 5. LO STALLO POLITICO** 29 maggio 2016
Dall'imperatore dello spread alle grivolette nel governo, cronaca di una giornata di arruolata

Motivation

<http://demo.istat.it/>

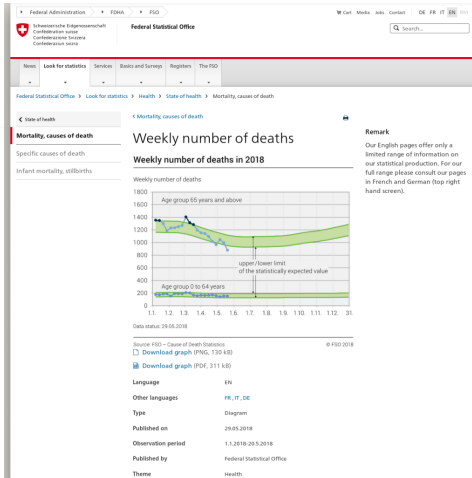


<https://github.com/gcampanella/istat-demographics>

The screenshot shows the ISTAT website interface with the following sections:

- popolazione residente**: Popolazione Residente per età, sesso e stato civile al 1° gennaio. Anno 2017, Anno 2016, Anno 2015, Anno 2014, Anno 2013, Anno 2012.
- bilancio demografico**: Bilancio Demografico e popolazione residente per sesso al 31 dicembre. Anno 2016, Anno 2015, Anno 2014, Anno 2013, Anno 2012, Anno 2011, Anno 2010, Anno 2009, Anno 2008, Anno 2007, Anno 2006, Anno 2005, Anno 2004, Anno 2003, Anno 2002, Anno 2001, Anno 2000, Anno 1999, Anno 1998, Anno 1997, Anno 1996, Anno 1995, Anno 1994, Anno 1993, Anno 1992, Anno 1991, Anno 1990, Anno 1989, Anno 1988, Anno 1987, Anno 1986, Anno 1985, Anno 1984, Anno 1983, Anno 1982, Anno 1981, Anno 1980, Anno 1979, Anno 1978, Anno 1977, Anno 1976, Anno 1975, Anno 1974, Anno 1973, Anno 1972, Anno 1971, Anno 1970, Anno 1969, Anno 1968, Anno 1967, Anno 1966, Anno 1965, Anno 1964, Anno 1963, Anno 1962, Anno 1961, Anno 1960, Anno 1959, Anno 1958, Anno 1957, Anno 1956, Anno 1955, Anno 1954, Anno 1953, Anno 1952, Anno 1951, Anno 1950.
- cittadini stranieri**: Popolazione straniera residente al 1° gennaio per età e sesso. Anno 2017, Anno 2016, Anno 2015, Anno 2014, Anno 2013, Anno 2012.
- dati precensuari**: I dati precensuari al 9 ottobre 2011 sono disponibili alla pagina [SCELE PRECENSUARIE DELLA POPOLAZIONE RESIDENTE NEI COMUNI \(2002-2011\)](#).
- elaborazioni**: Tavole di Mortalità della popolazione per provincia e regione di residenza. Anno 1974-2016. Previsioni della popolazione. Anno 2017-2065. Bilanciamento della popolazione per età e sesso al 1° gennaio. Anno 2002-2011. Anno 1992-2001. Anno 1982-1991. Bilanciamento del bilancio demografico per sesso. Anno 2001-2011. Anno 1991-2001. Tavole di Fecundità della popolazione totale per regione di residenza. Anno 1952-2004.
- dati Cuni**: Espatriati e divieti di espatrio. La rilevazione delle espatriazioni e dei divieti - Anno 2009-2012. Iscritti in anagrafe per nascita. Dati relativi agli anni 1990-2016. Iscritti e cancellati all'anagrafe per trasferimenti di residenza. Dati relativi agli anni 2001-2016. Cancellati dall'anagrafe per decesso. Anno 2011-2016. I sottratti. Dati relativi agli anni 2004-2016. Cittadini non censurati regolarmente presenti in Italia. Dati relativi agli anni 2008-2011. Permessi di soggiorno. Permessi di soggiorno al 1° gennaio. Anno 1992-2011. Bilanci demografici dei cittadini stranieri. Dati al 1° gennaio 2009. Le nascite in Italia. La rilevazione delle nascite di tutte le donne. Anno 1991 e 1996. Indagine componente sulle nascite. I risultati dell'indagine demografica sulle nascite degli anni 2002 e 2005. Indicatori demografici.

Motivation



Original data

- Births, deaths, and net migration
- Monthly resolution from January 2004 till November 2017
- At municipality (*comune*) level
- Stratified by sex

Aggregated data

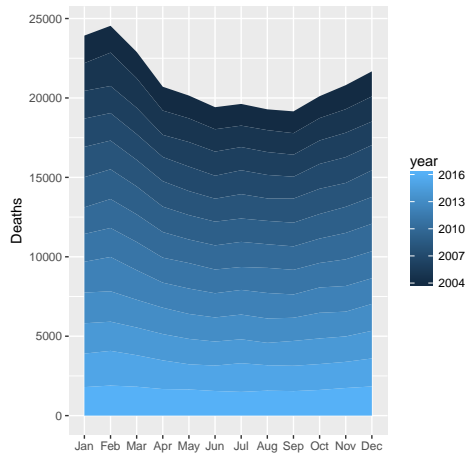
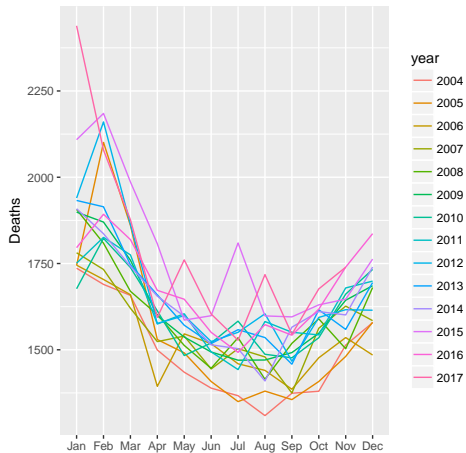
- Deaths only
- Monthly resolution from January 2004 till November 2017
- At **region** level ($N = 20$)
- Stratified by sex

| | Start | End | Length |
|-----------------|--------------|---------------|---------------|
| Training | January 2004 | June 2016 | 12.5 years |
| Test | July 2016 | November 2017 | 17 months |

Data are **unnormalised** monthly counts

- Boundary changes
- Population size (pre-census vs post-census)
- Calendar adjustment

Exploratory data analysis



Models

| Family | Method | Package |
|--------------|------------------|----------|
| Baseline | Naïve (RW) | forecast |
| | Seasonal naïve | forecast |
| | Naïve with drift | forecast |
| | Average | forecast |
| Univariate | ETS | forecast |
| | ARIMA | forecast |
| | BSTS | bsts |
| | Prophet | prophet |
| Hierarchical | HTS | hts |

Modelling

Naïve and average methods

For all $h = 1, 2, \dots$,

Naïve (RW)

$$\hat{y}_{T+h|T} = y_T$$

Seasonal naïve with period m

$$\hat{y}_{T+h|T} = y_{T+h-m(\lfloor (h-1)/m \rfloor + 1)}$$

Naïve with drift

$$\hat{y}_{T+h|T} = y_T + h(y_T - y_1)/(T - 1)$$

Average

$$\hat{y}_{T+h|T} = \sum_{t=1}^T y_t / T$$

Time series decomposition

Common components

- Trend-cycle T_t
- Seasonal S_t
- Remainder R_t

Additive model

$$y_t = T_t + S_t + R_t$$

Multiplicative model

$$y_t = T_t \times S_t \times R_t$$

Modelling

Exponential smoothing

Simple exponential smoothing (SES)

Given a smoothing parameter $0 \leq \alpha \leq 1$,

$$\hat{y}_{t+1|t} = \alpha y_t + (1 - \alpha) \hat{y}_{t|t-1}$$

$$\hat{y}_{t+h|t} = l_t \quad \text{(forecast)}$$

$$l_t = \alpha y_t + (1 - \alpha) l_{t-1} \quad \text{(smoothing)}$$

Holt's linear trend method

Given a smoothing parameter $0 \leq \beta \leq 1$,

$$\begin{aligned}\hat{y}_{t+h|t} &= l_t + hb_t && \text{(forecast)} \\ l_t &= \alpha y_t + (1 - \alpha)(l_{t-1} + b_{t-1}) && \text{(level)} \\ b_t &= \beta(l_t - l_{t-1}) + (1 - \beta)b_{t-1} && \text{(trend)}\end{aligned}$$

Gardner and McKenzie's damped trend method

Given a damping parameter $0 < \phi < 1$,

$$\hat{y}_{t+h|t} = l_t + (\phi + \phi^2 + \dots + \phi^h)b_t \quad \text{(forecast)}$$
$$l_t = \alpha y_t + (1 - \alpha)(l_{t-1} + \phi b_{t-1}) \quad \text{(level)}$$
$$b_t = \beta(l_t - l_{t-1}) + (1 - \beta)\phi b_{t-1} \quad \text{(trend)}$$

Holt-Winters' seasonal (additive) method

Given a smoothing parameter $0 \leq \gamma \leq 1$ and a frequency $m \in \mathbb{N}$,

$$\hat{y}_{t+h|t} = l_t + hb_t + s_{t+h-m(\lfloor (h-1)/m \rfloor + 1)} \quad (\text{forecast})$$

$$l_t = \alpha(y_t - s_{t-m} + (1 - \alpha)(l_{t-1} + b_{t-1})) \quad (\text{level})$$

$$b_t = \beta(l_t - l_{t-1}) + (1 - \beta)b_{t-1} \quad (\text{trend})$$

$$s_t = \gamma(y_t - l_t) + (1 - \gamma)s_{t-m} \quad (\text{seasonality})$$

ETS methods

- **Error**
 - Additive
 - Multiplicative
- **Trend**
 - None
 - Additive
 - Additive damped
- **Seasonality**
 - None
 - Additive
 - Multiplicative

⇒

$2 \times 3 \times 3 = 18$
possible configurations

Modelling

ARIMA models

Backshift operator \mathcal{B}

Let's introduce the **backshift operator** \mathcal{B} ,

$$\mathcal{B}y_t = y_{t-1}$$

$$\mathcal{B}^2 y_t = y_{t-2}$$

$$\vdots$$

$$\mathcal{B}^m y_t = y_{t-m}$$

Backshift operator \mathcal{B}

We can rewrite first-order differences in terms of \mathcal{B} ,

$$\begin{aligned}y_t - y_{t-1} &= y_t - \mathcal{B}y_t \\ &= (1 - \mathcal{B})y_t\end{aligned}$$

In general, \mathcal{B} follows algebraic rules,

$$\begin{aligned}(1 - \mathcal{B})(1 - \mathcal{B}^m)y_t &= (1 - \mathcal{B}^m - \mathcal{B} + \mathcal{B}^{m+1})y_t \\ &= y_t - y_{t-m} - y_{t-1} + y_{t-m-1} \\ &= (y_t - y_{t-m}) - (y_{t-1} - y_{(t-1)-m})\end{aligned}$$

Autoregressive and moving average models

Autoregressive AR(p) model of order p

$$y_t = \beta_0 + \beta_1 y_{t-1} + \dots + \beta_p y_{t-p} + \epsilon_t$$

Moving average MA(q) model of order q

$$y_t = \gamma_0 + \gamma_1 \epsilon_{t-1} + \dots + \gamma_q \epsilon_{t-q} + \epsilon_t$$

Non-seasonal ARIMA(p, d, q) model

$$(1 - \beta_1\mathcal{B} - \dots - \beta_p\mathcal{B}^p)(1 - \mathcal{B})^d y_t = \alpha + (1 + \gamma_1\mathcal{B} + \dots + \gamma_q\mathcal{B}^q)\epsilon_t$$

Non-seasonal ARIMA(p, d, q) model

$$(1 - \beta_1\mathcal{B} - \dots - \beta_p\mathcal{B}^p)(1 - \mathcal{B})^d y_t = \alpha + (1 + \gamma_1\mathcal{B} + \dots + \gamma_q\mathcal{B}^q)\epsilon_t$$

Seasonal ARIMA(p, d, q)(P, D, Q) $_m$ model

$$(1 - \beta_1\mathcal{B} - \dots - \beta_p\mathcal{B}^p)(1 - B_1\mathcal{B}^m - \dots - B_P\mathcal{B}^{Pm})(1 - \mathcal{B})^d(1 - \mathcal{B}^D) y_t \\ = \alpha + (1 + \gamma_1\mathcal{B} + \dots + \gamma_q\mathcal{B}^q)(1 + \Gamma_1\mathcal{B}^m + \dots + \Gamma_Q\mathcal{B}^{Qm})\epsilon_t$$

Modelling

Other methods

Bayesian Structural Time Series (BSTS) models

- Introduced by S. L. Scott and H. Varian (Google)
- Ensemble method
- Structural time series model + regression component

Model evaluated

- Local linear trend
- Seasonal model with $m = 12$

Prophet

- Introduced by S. J. Taylor and B. Letham (Facebook)
- Curve fitting (similarly to GAMs)
- Decomposition into trend, seasonality, and holidays

Model evaluated

- Default settings
- No daily or weekly seasonality

Hierarchical time series models

- Introduced by R. J. Hyndman et al. (Monash University)
- Independent forecasts + aggregation at different levels
- Many different aggregation methods

Models evaluated

- Forecasting methods: ARIMA, ETS, RW
- 5 aggregation methods \times 4 weighting schemes

Modelling

Measures

Scale-dependent measures

Given the prediction errors $e_{T+h} = y_{T+h} - \hat{y}_{T+h}$, ...

Measure

Mean absolute error $\text{mean}(|e_t|)$

Root-mean-square error $\sqrt{\text{mean}(e_t^2)}$

Percentage errors

Given the **percentage** errors $p_t = 100e_t/y_t, \dots$

Measure

Mean absolute percentage error $\text{mean}(|p_t|)$

Symmetric MAPE $\text{mean}(200|y_t - \hat{y}_t|/(y_t + \hat{y}_t))$

Scaled errors

Given the **scaled** errors...

$$q_t = \frac{e_t}{\frac{1}{T-1} \sum_{t'=2}^T |y^{t'} - y^{t'-1}|} \quad \text{or} \quad q_t = \frac{e_t}{\frac{1}{T-m} \sum_{t'=m+1}^T |y^{t'} - y^{t'-m}|},$$

the **mean absolute scaled error** is simply $\text{mean}(|q_t|)$

Interpretation

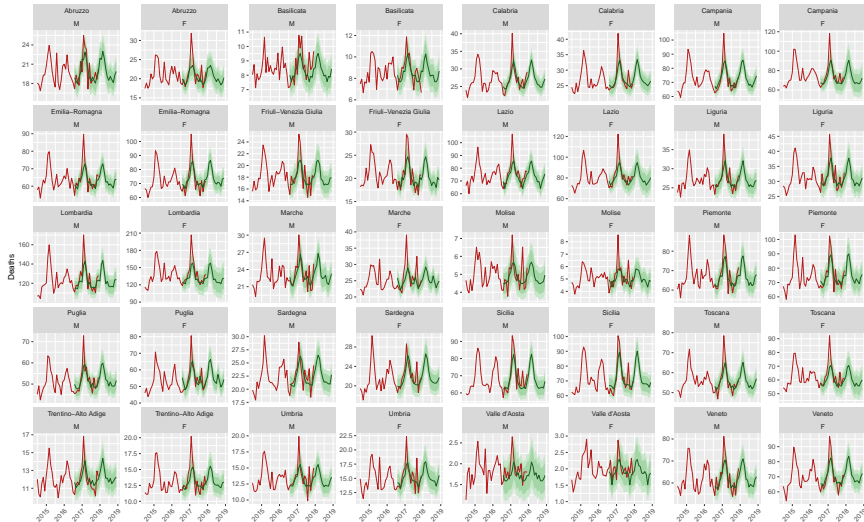
For $q_t < 1$, the forecast is better than the average (seasonal) naïve forecast (computed on the training data)

Results

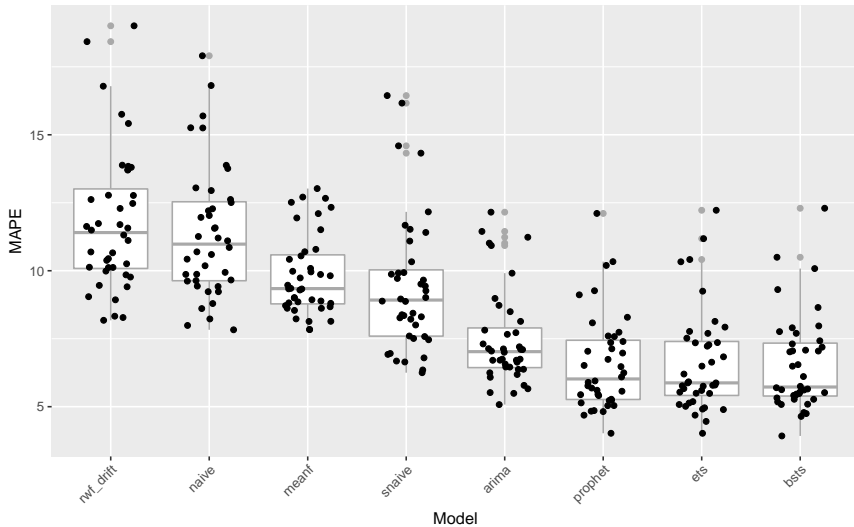
Seasonal naïve forecasts



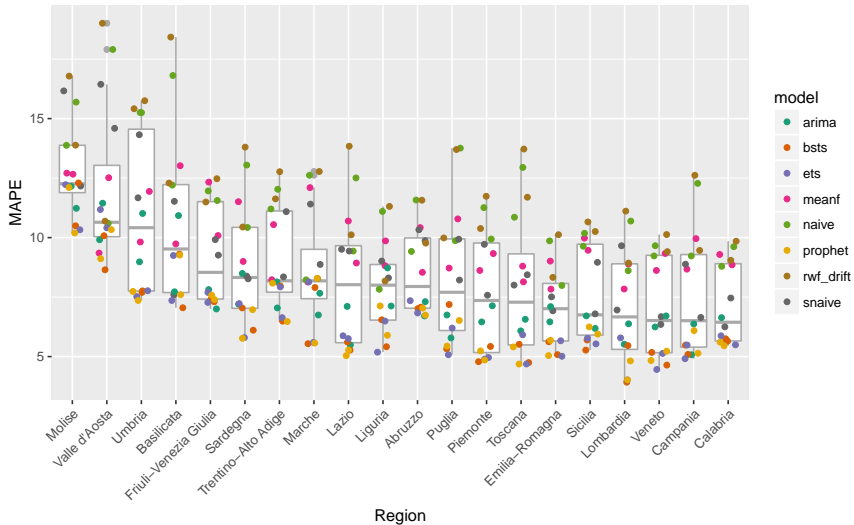
ETS forecasts



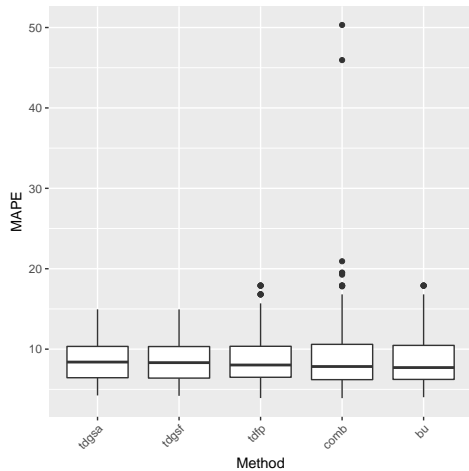
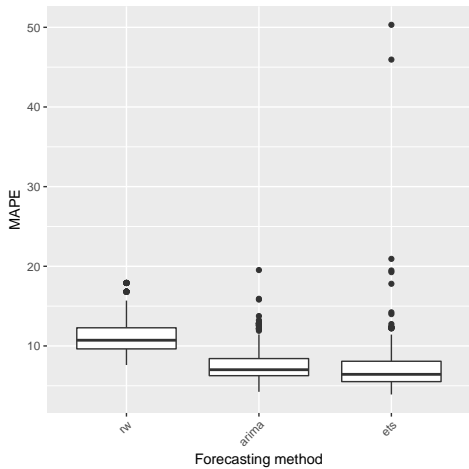
Univariate models



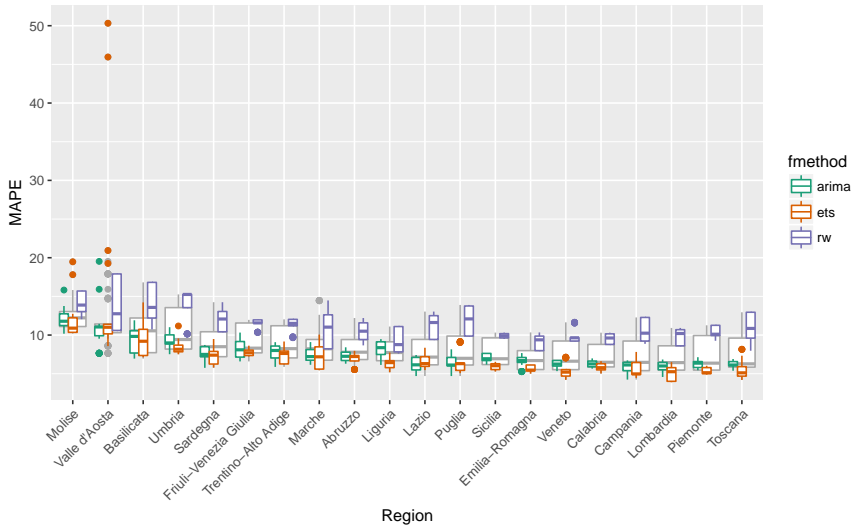
Univariate models



HTS models



HTS models



And the winner is...

| Method | MAPE |
|---------------------|-------|
| BSTS | 6.52% |
| Prophet | 6.58% |
| ETS | 6.62% |
| HTS (bottom-up ETS) | 6.62% |
| ARIMA | 7.49% |
| Seasonal naïve | 9.44% |
| Average | 9.83% |
| Naïve (RW) | 11.4% |
| Naïve with drift | 11.8% |

Recommendations

Time series are messy!

- Temporal resolution and spacing
- Calendar adjustment
- Model evaluation and cross-validation
- Hierarchical structure

Time series are fun!

- Data visualisation
- Models (often) interpretable
- Anomaly detection

Before you get started...

Ask yourself...

- Do I have **enough data**?

Before you get started...

Ask yourself...

- Do I have enough data?
- Is my time series **evenly spaced**?

Before you get started...

Ask yourself...

- Do I have enough data?
- Is my time series evenly spaced?
- Which **measures** do I care about?

During modelling...

- **Visualise** — Trend? Seasonality? ‘Spikes’?

During modelling...

- Visualise — Trend? Seasonality? ‘Spikes’?
- Start with a **simple** model

During modelling...

- Visualise — Trend? Seasonality? ‘Spikes’?
- Start with a simple model
- Plot the ACF of residuals

Future work

- Compare even more models (including neural networks)
- Include exogenous covariates such as temperature
- Build a user interface