

Basics of Time-Series Analysis

Charalampos E. Tsourakakis
ctsourak@bu.edu

CS365 Foundations of Data Science

April 2024

Time series – ubiquitous!

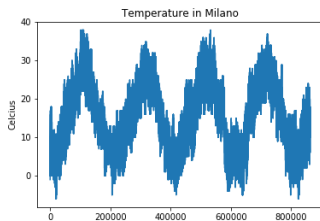
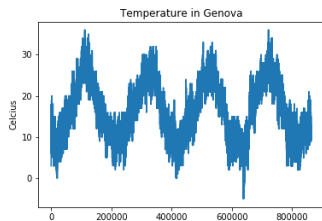
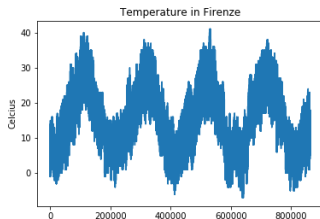
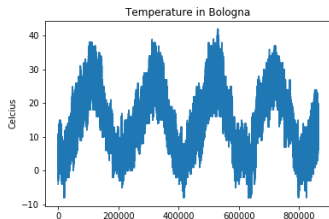
Time-series appear in any domain that involves temporal measurements.

Examples of such domains include:

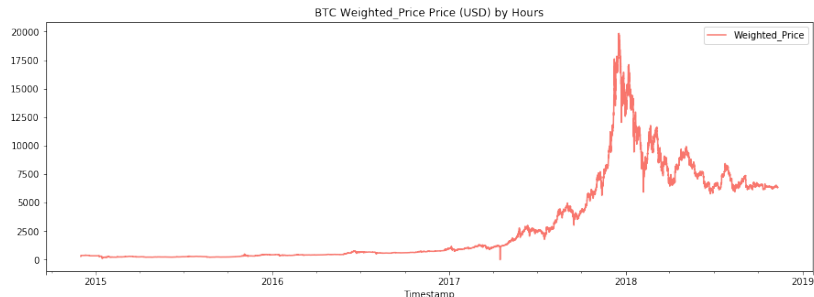
- signal processing
- econometrics
- mathematical finance
- weather forecasting
- electroencephalography
- control engineering
- astronomy
- communications engineering
- ...

Time series – examples

Temperature time series across 4 Italian cities



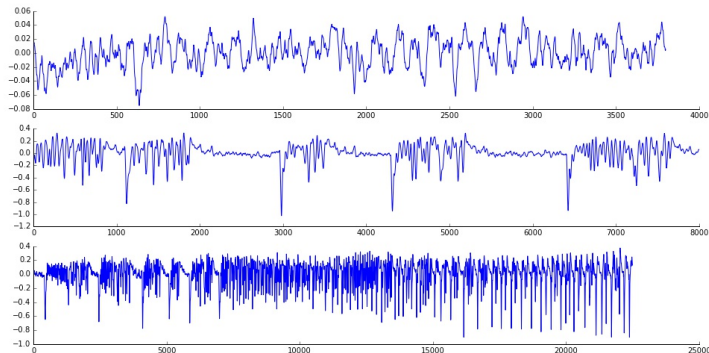
Time series – examples



Bitcoin price

- When we have a *single time-series* we wish to forecast using its past values, we call this *univariate* analysis.

Time series – examples



EEG time series from different sections of the brain

- When we have a *target time-series* we wish to forecast using its past values and other co-evolving time-series, we call this *multivariate analysis*.

Time series analysis

Why do we analyze time-series?

① Interpretation

E.g., seasonal adjustment

② Forecasting

E.g., predict stock prices

③ Control

E.g., how does increasing VAT will affect (un)employment?

Time series analysis

④ Hypothesis testing

E.g., should we be believers in global warming?



⑤ Understanding catastrophic events

E.g., will an impactful earthquake take place tomorrow (here)?

Today's focus – Forecasting



*A trend is a trend is a trend. But the question is, will it bend?
Will it alter its course through some unforeseen force and
come to a premature end? (Alex Cairncross)*

Today's outline

① Part I: Fundamentals

- Stochastic processes, strict vs weak stationarity, correlograms, partial correlograms, periodograms, autoregressive models (AR), moving average models (MA), transformations, compact time-series description.

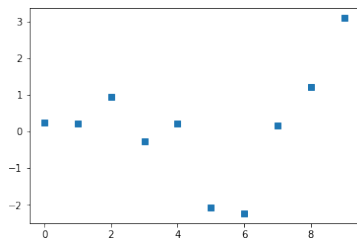
② Part II: Forecasting methods

- Evaluation protocols, basic metrics, null models, spectral methods, ARMA, ARIMA, SARIMA, VAR, similarity search, deep learning (deep feedforward neural nets, recursive neural nets)

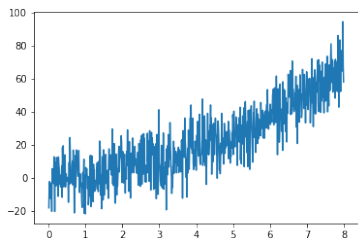
Stochastic processes – informal description

Definition. A stochastic process is a *collection of random variables indexed by time*.

- (a) Discrete time stochastic processes (e.g., X_0, X_1, X_2, \dots)
- (b) Continuous time stochastic processes ($\{X_t\}_{t \geq 0}$)



(a)



(b)

Time-series – formal definition

- A **time series** is a stochastic process consisting of random variables indexed by time t .
- The stochastic behavior of the *time series* is determined by specifying the probability density/mass functions (pdf's)

$$p(x_{t_1}, x_{t_2}, \dots, x_{t_m}),$$

for all finite collections of time indexes

$$\{(t_1, \dots, t_m), m < +\infty\},$$

i.e, all finite-dimensional distributions of $\{X_t\}$.

Time-series – Stationarity

Strictly stationary time-series. A time series $\{X_t\}$ is strictly stationary if $\forall t, m, (t_1, \dots, t_m)$

$$p(t_1 + \tau, \dots, t_m + \tau) = p(t_1, \dots, t_m).$$

- In other words, a time series is strictly stationary if the probability distribution is invariant under time translation.

Examples

- (a) *iid* processes are strictly stationary.
- (b) $X_t = Z_1 \cos(t) + Z_2 \sin(t)$ is strictly stationary if Z_1, Z_2 are independent normal variables.
- (c) Random walks in certain types of graphs (stationary Markov chains)

Remark: Stationary time series are typically non-stationary.

Time-series – Covariance stationarity

Covariance stationary time-series. A time series $\{X_t\}$ is covariance stationary if

$$\mathbb{E}[X_t] = \mu$$

$$\text{Var}[X_t] = \sigma^2$$

$$\text{Cov}[X_t, X_{t+h}] = \gamma(h)$$

Reminder: $\text{Cov}[X_t, X_{t+h}] = \mathbb{E}[(X_t - \mu_t)(X_{t+h} - \mu_{t+h})]$.

Basic exploration tools – auto-correlation function

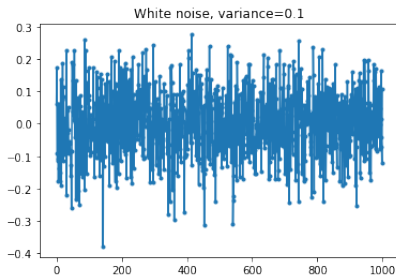
Auto-correlation function of $\{X_t\}$ is defined as

$$\begin{aligned}\rho_X(h) &= \frac{\gamma_X(h)}{\gamma_X(0)} \\ &= \frac{\text{Cov}(X_{t+h}, X_t)}{\text{Cov}(X_t, X_t)} \\ &= \frac{\text{Cov}(X_{t+h}, X_t)}{\text{Var}[X_t]} \\ &= \text{Corr}(X_{t+h}, X_t)\end{aligned}$$

Basic exploration tools – auto-correlation function

White noise: $X_t \sim WN(0, \sigma^2)$.

- $\mathbb{E}[X_t] = 0$
- $\text{Var}[X_t] = \sigma^2$
- $\Pr[X_t \leq x_t] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x_t} e^{-x^2/2} dx$.



Question: What is the auto-correlation function of white noise? Is it stationary?

Basic exploration tools – auto-correlation function

We have:

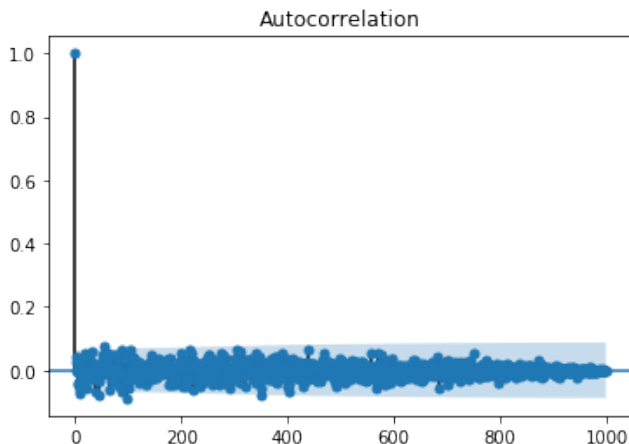
$$\gamma_X(t+h, t) = \begin{cases} \sigma^2 & h = 0 \\ 0 & h > 0 \end{cases}$$

Thus,

- $\mu_t = 0$ for all t
- $\gamma_X(t+h, t) = \gamma_X(h, 0)$ for all $t \geq 0$
- Thus $\rho_X(h) = 1$ if $h = 0$, 0 otherwise.
- **question:** is X_t stationary?

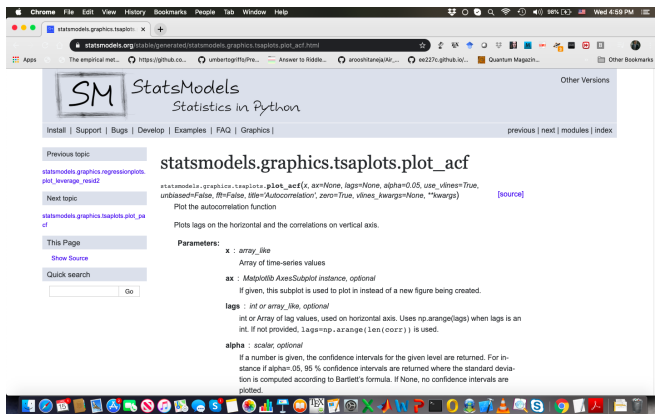
Basic exploration tools – correlogram

- The **correlogram** plots the auto-correlation function versus the lag.



Basic exploration tools – correlogram

Python:



The screenshot shows a web browser displaying the StatsModels website. The page title is "StatsModels Statistics in Python". The main content area is titled "statsmodels.graphics.tsaplots.plot_acf". The page includes a navigation menu with links for "Install", "Support", "Bugs", "Develop", "Examples", "FAQ", and "Graphics". The main content area contains the following text:

statsmodels.graphics.tsaplots.plot_acf(x, ax=None, lags=None, alpha=0.05, use_vlines=True, unbiased=False, fit=False, title='Autocorrelation', zero=True, vlines_kwargs=None, **kwargs) [source]

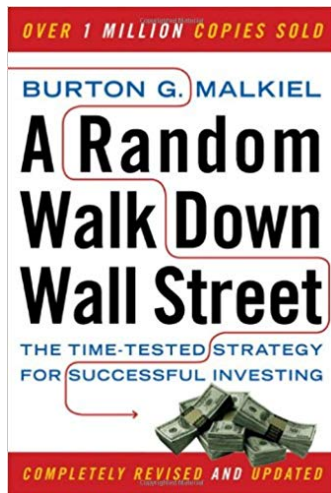
Plot the autocorrelation function

Plots lags on the horizontal and the correlations on vertical axis.

Parameters:

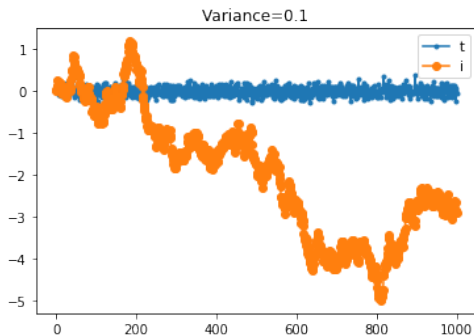
- x** : array_like
Array of time-series values
- ax** : Matplotlib AxesSubplot instance, optional
If given, this subplot is used to plot in instead of a new figure being created.
- lags** : int or array_like, optional
int or Array of lag values, used on horizontal axis. Uses np.arange(lags) when lags is an int. If not provided, lags=np.arange(1+len(xcorr)) is used.
- alpha** : scalar, optional
If a number is given, the confidence intervals for the given level are returned. For instance if alpha=.05, 95 % confidence intervals are returned where the standard deviation is computed according to Bartlett's formula. If None, no confidence intervals are plotted.

What is a random walk time series?



Random walk

Random walk: $S_t = \sum_{i=1}^t X_i$ for $X_i \sim WN(0, \sigma^2)$.



question: is a random walk stationary or not?

Random walk is not stationary

- $\mathbb{E}[S_t] = 0$
- However, the covariance function is a function of t :

$$\begin{aligned}\gamma_S(t+h, t) &= \text{Cov}(S_{t+h}, S_t) \\ &= \text{Cov}\left(S_t + \sum_{s=1}^h X_{t+s}, S_t\right) \\ &= \text{Cov}(S_t, S_t) = \text{Var}[S_t] = t\sigma^2.\end{aligned}$$

Therefore, S_t is not stationary.

Moving average process $MA(1)$

We define the moving average process of order 1 $MA(1)$ as

$$X_t = Z_t + \theta Z_{t-1}, \{Z(t)\} \sim WN(0, \sigma^2).$$

- We have $\mathbb{E}[X_t] = 0$ and

$$\begin{aligned}\gamma_X(t+h, t) &= \mathbb{E}[X_{t+h}X_t] \\ &= \mathbb{E}[(Z_{t+h} + \theta Z_{t+h-1})(Z_t + \theta Z_{t-1})] \rightarrow\end{aligned}$$

$$\gamma_X(t+h, t) = \begin{cases} \sigma^2(1 + \theta^2) & \text{if } h = 0 \\ \sigma^2\theta & \text{if } h = 1 \\ 0 & \text{otherwise} \end{cases}$$

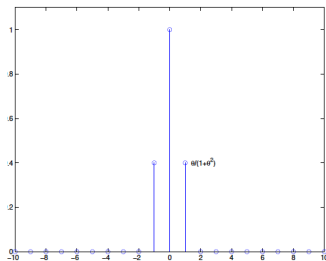
Moving average process $MA(1)$

- As we observe $MA(1)$ is stationary.
 - $\mathbb{E}[X_t] = 0$
 - $\gamma_X(t+h, t) = \gamma_X(h, 0)$ is
- How *can we use it?*
 - $\Delta \text{Icecream}_t = \Delta \text{temperature}_t - 0.9 \Delta \text{temperature}_{t-1}$,
- We assume that the changes in the temperature are iid normal random variables.

- **Important observation:** We model differences, not the actual time series values. This is a common time-series transformation that we use frequently in practice.

Moving average process $MA(1)$

The correlogram of ρ_X of an $MA(1)$ process:



- Moving average processes generalize to order q $MA(q)$ processes

$$X_t = \mu + Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}.$$

- The autocorrelation function of an $MA(q)$ process is zero at lag $q + 1$ and greater.

Autoregressive process $AR(1)$

Mathematically an $AR(1)$ process is defined as

$$X_t = \rho X_{t-1} + Z_t, Z_t \sim WN(0, \sigma^2).$$

- In general, the order p autoregressive process is defined similarly as:

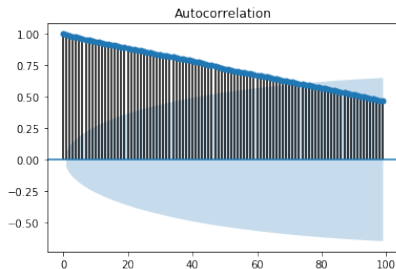
$$X_t = c + \rho_1 X_{t-1} + \dots + \rho_p X_{t-p} + Z_t, Z_t \sim WN(0, \sigma^2).$$

- **Exercise**

① $\mathbb{E}[X_t] = \frac{c}{1-\phi}$

② $\gamma_X(t+h, t) = \frac{\sigma^2}{1-\phi^2} \phi^h.$

Autoregressive process $AR(1)$



- In contrast to moving average processes, the ACF plot cannot tell us (at least by inspection) the order p (here, $p = 1$).
- Still though, we will show another tool, the partial correlogram plot that allows us to get an idea of the order of the AR process.

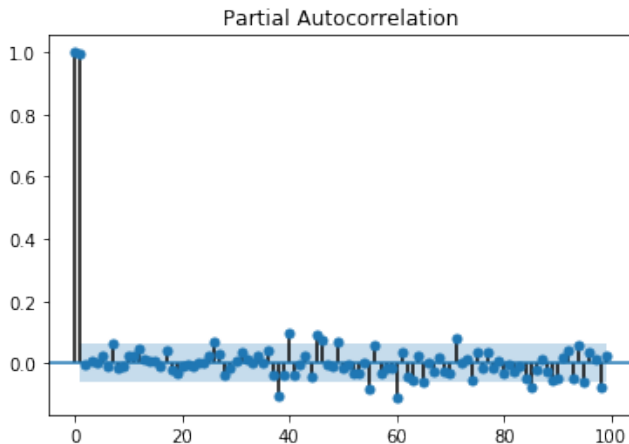
Basic exploration tools - partial correlogram

- To understand the necessity of another tool that is used *jointly* with the correlogram we will give an example.

$$x(t+2) = x(t+1) + \epsilon(t+1) = x(t) + \epsilon(t) + \epsilon(t+1).$$

- Thus $x(t+2)$ and $x(t)$ are only related because of $x(t+1)$ in between.
- The idea of the partial correlogram is to measure the correlation of $x(t+h)$ and $x(t)$ after removing linear relationship due to values in between.

Partial correlogram $AR(1)$



Estimating the ACF

- In reality, we have only access to the data.
- Estimating the mean:

$$\bar{x} = \frac{1}{n} \sum_{t=1}^n x_t.$$

- Estimating autocovariance function:

$$\hat{\gamma}(h) = \frac{1}{n} \sum_{t=1}^{n-h} (x_{t+h} - \bar{x}) \cdot (x_t - \bar{x}).$$

- Estimating autocorrelation function:

$$\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)}.$$

Stationarity statistical tests

- Given the data, we can test whether the time series (or the differences, or other transformations) is stationary or not.
 - Dickey-Fuller test (unit root test)
 - Kwiatkowski-Phillips-Schmidt-Shin test
- The KPSS test can test *trend-stationarity*
- Dickey-Fuller and its variations are popular.
 - The Dickey-Fuller test:
 - formally, it tests the null hypothesis that a unit root is present in an autoregressive model.
 - Alternative hypothesis: time-series is stationary

Dickey-Fuller test in Python

```
1 from statsmodels.tsa.stattools import adfuller
2
3
4 def Dickey_Fuller_test(timeseries):
5     '''
6     Test stationarity
7     '''
8     #Perform Dickey-Fuller test:
9     print( 'Results of Dickey-Fuller Test:')
10    dftest = sm.tsa.stattools.adfuller(
timeseries)
11    dfoutput = pd.Series(dftest[0:4], index=['
Test Statistic', 'p-value', '#Lags Used', '
Number of Observations Used'])
12    for key,value in dftest[4].items():
13        dfoutput['Critical Value (%s)'%key] =
value
14    print(dfoutput)
```

How to interpret the Dickey-Fuller test?

- Null Hypothesis: time series has a unit root, meaning it is non-stationary.
- Alternative Hypothesis : time-series is stationary
- Code returns:

```
1 Results of Dickey-Fuller Test:
2 Test Statistic           -1.102293e+01
3 p-value                   5.925157e-20
4 Number of Observations Used  9.900000e+01
5 Critical Value (5%)       -2.891208e+00
6 Critical Value (1%)       -3.498198e+00
7 Critical Value (10%)      -2.582596e+00
```

- Check the value of the statistic and compare it the critical values.
- If it is less than the critical value, then we can reject the null hypothesis with that level of confidence.
- Also p -values are informative: Rule of thumb:
 - If $p\text{-value} > 0.05$, then time-series is non-stationary.

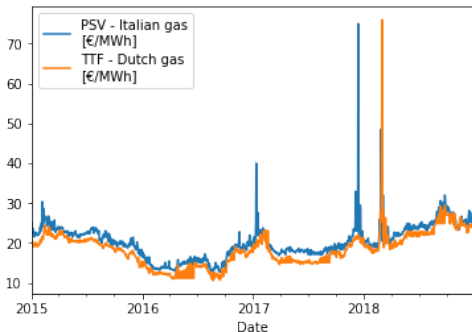
Kwiatkowski-Phillips-Schmidt-Shin test in Python

```
1 def kpss_test(timeseries):
2     print ('Results of KPSS Test:')
3     kpsstest = tsa.stattools.kpss(timeseries)
4     kpss_output = pd.Series(kpsstest[0:2], index
5                             =['Test Statistic', 'p-value'])
6     for key,value in kpsstest[2].items():
7         kpss_output['Critical Value (%s)'%key] =
            value
8     print (kpss_output)
```

Cointegration – Stationarity in multivariate time series

- When we have multiple co-evolving time-series, there is the notion of cointegration.
- For example:
 - D_t : position of a dog at time t
 - P_t : position of the person who takes the dog for a walk at time t
- While D_t, P_t may look individually random, their difference is not.
- **Question:** does there exist a constant c such that $P_t - cD_t$ is stationary?

Cointegration – Stationarity in multivariate time series

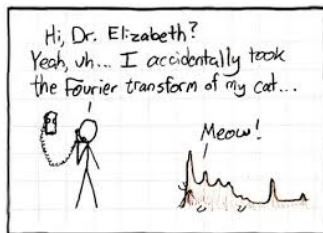


```
1 from statsmodels.tsa.stattools import coint
2 coint(P,Q)
```

Fourier analysis - basics



Joseph Fourier showed how to represent a periodic function as a sum of trigonometric functions (oscillations).



Fourier analysis - basics

$$x(t) - x(\bar{t}) = \sum_{k=1}^n \gamma_k \left(a_k \cos(2\pi t\omega_k) + b_k \sin(2\pi t\omega_k) \right).$$

where $\omega_k = \frac{k}{n}$. If we define:

- $C_k = \sqrt{a_k^2 + b_k^2}$
- $\phi_k = \arctan \frac{b_k}{a_k}$

then we can also write the formula also as (trigonometric identity):

$$x(t) - x(\bar{t}) = \sum_k \gamma_k C_k \cos(2\pi t\omega_k - \phi_k).$$

Fourier analysis - sinusoidal decomposition

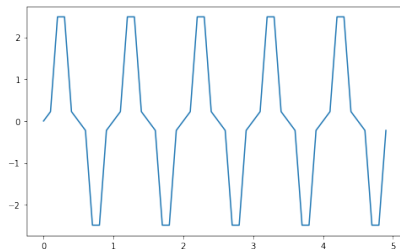
- We prefer to express things in *complex numbers*.
- Trigonometric identities are trivial. E.g.,

$$\begin{aligned}\sin(x) \cos(x) &= \frac{1}{2j}(e^{ix} - e^{-ix}) \frac{1}{2}(e^{ix} + e^{-ix}) = \\ &= \frac{1}{4j}(e^{i2x} - 1 + 1 - e^{-2ix}) \\ &= \frac{1}{2} \frac{1}{2j}(e^{i2x} - e^{-2ix}) = \frac{1}{2} \sin(2x).\end{aligned}$$

Periodogram

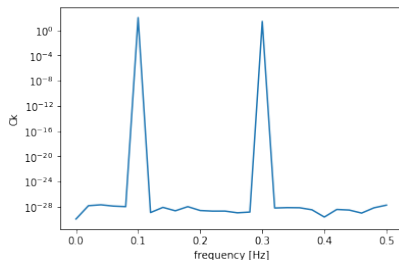
- A plot of nC_k^2 versus the frequency ω_k for $k = 1, 2, \dots$ is called the **periodogram** of the dataset.
- The **periodogram** is a very useful tool.
- Random-signals that lack structure have all sinusoids with equal importance.
- If a time series has a strong sinusoidal signal for some frequency, then there will be a peak in the periodogram at that frequency.
- If a large C_k value appears at ω_k and other large values $C_{\ell k}$ appear at multiples $\ell\omega_k$ where $\ell = 2, \dots$ it suggests that there is a non-sinusoidal signal with that specific frequency.

Periodogram – example



```
1 sig = np.sin(2 * np.pi * f1*time_vec ) + 2* np
    .sin(2 * np.pi * f2*time_vec )
2 f, Pxx_den = signal.periodogram(sig)
3 plt.semilogy(f, Pxx_den)
4 plt.show()
```


Periodogram



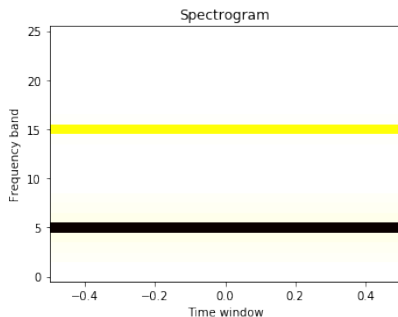
- Python plots

$$\hat{f}(\omega) = \frac{1}{n} \left| \sum_{t=1}^n x(t) e^{2\pi i t \omega} \right|, \omega \in [0, 0.5].$$

- Symmetry $\hat{f}(\omega) = \hat{f}(1 - \omega)$

Spectrogram

- **Which frequencies appear across time?**
- Spectrogram!

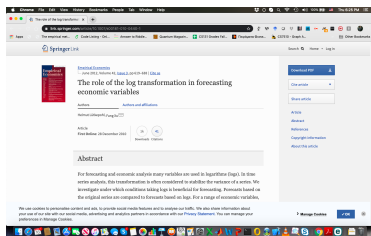


Spectrogram – Python snippet

```
1 from scipy import signal
2 freqs, times, spectrogram = signal.spectrogram(
    sig)
3
4 plt.figure(figsize=(5, 4))
5 plt.imshow(spectrogram, aspect='auto', cmap='
    hot_r', origin='lower')
6 plt.title('Spectrogram')
7 plt.ylabel('Frequency band')
8 plt.xlabel('Time window')
9 plt.tight_layout()
```

Time-series transformations

- Earlier we saw that a random walk time series is not stationary. Nonetheless, the increments form a stationary time series (white noise).
- Therefore, *transforming* the time-series **can** give a different structured signal
- **Goal:** Transform the signal so that hopefully it is easier to model mathematically (and hence better forecasts).



Time-series transformations - Box-Cox

- A Box-Cox transformation is a way to stabilize the variance of a time series with non-negative values.
- If negative values are present, we use the Yeo-Johnson power transform.
- Box-Cox function is invertible.
- Working with the Box-Cox transformation of financial time series frequently helps (even slightly).
- There is a parameter λ that defines the Box-Cox transformation of a measurement y as follows:

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

Remark

- An example of how much a transformation can help the prediction on the gas prices in Italy.

```
1 Box-Cox transformation
2 -----
3 Absolute error:12.727407888799966
4 -----
5
6 No preprocessing
7 -----
8 Absolute error:23.151990039645284
9 -----
```

Time-series transformations - Box-Cox

```
1 def box_cox(y):
2     if not isinstance(y, pd.Series):
3         y = pd.Series(y)
4         y.astype(float)
5     y_boxcox, lambda = stats.boxcox(y)
6     return y_boxcox, lambda
7
8
9 def invboxcox(y, lambda):
10    if not isinstance(y, pd.Series):
11        y = pd.Series(y)
12        y.astype(float)
13
14    if lambda == 0:
15        return(np.exp(y))
16    else:
17        return(np.exp(np.log(lambda*y+1)/lambda))
```

Time-series transformations - Differencing

- **First order differencing:** convert time series $x(1), x(2), \dots, x(n)$ to $x(2) - x(1), \dots, x(n) - x(n - 1)$.

```
1 def first_order_diff(y):
2     if not isinstance(y, pd.Series):
3         y = pd.Series(y)
4         y.astype(float)
5     z = y.diff(1).dropna()
6     return z
7
8
9
10 def first_order_diff_numpy(data):
11     return [data[i] - data[i - 1] for i in range
12             (1, len(data))]
```


Time-series transformations - Differencing (variations)

- **First order differencing:** sometimes when we assume increments have seasonality we may do the following:
E.g., weekly seasonality on daily data, then we convert time series $x(1), x(2), \dots, x(n)$ to $x(7) - x(1), x(8) - x(2), \dots, x(n) - x(n - 7)$

```
1 def difference(data, gap):  
2     return [data[i] - data[i - gap] for i in range  
            (gap, len(data))]
```

- Another variation: second order differencing, i.e., differences on the differences (third order etc.)

Time-series transformations - Logarithm

- The log-transformation of a time-series under certain mathematical modeling assumptions provably improves the accuracy of predictors.
- In practice, where the distribution that generates the time-series is not known, we use it, and evaluate it empirically.

```
1 def log_(y):  
2     ''' The logarithm of the time series (always  
3     apply on non-negative values!)'''  
4     if not isinstance(y, pd.Series):  
5         y = pd.Series(y)  
6         y.astype(float)  
7     log_y = np.log(y)  
8     return log_y
```

Time-series transformations - Log-returns

- Differences on the log results in what is known as log-returns.
- Transforms series (x_1, \dots, x_n) to $\log(x_2/x_1), \dots, \log(x_n/x_{n-1})$

```
1 def log_returns(y):  
2     if not isinstance(y, pd.Series):  
3         y = pd.Series(y)  
4         y.astype(float)  
5     log_returns = np.log(y/y.shift(1)).dropna()  
6     return log_returns  
7
```

Time-series transformations - Standardization

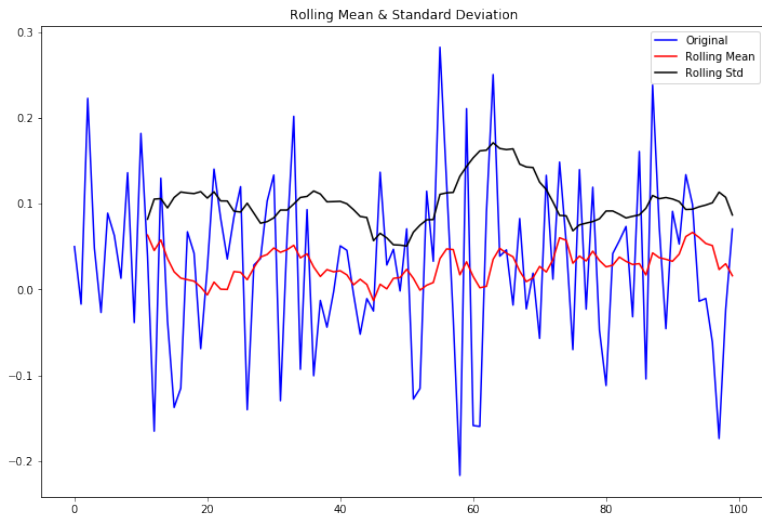
- Standardization is especially important when we have multiple time series, with different types of measurements.
- There exist two popular standardization methods:

$$x'_i = \frac{x_i - \min_{1 \leq j \leq n}(x_j)}{\max_{1 \leq j \leq n}(x) - \min_{1 \leq j \leq n}(x)},$$

and

$$x'_i = \frac{x_i - \bar{x}}{\sigma(x)},$$

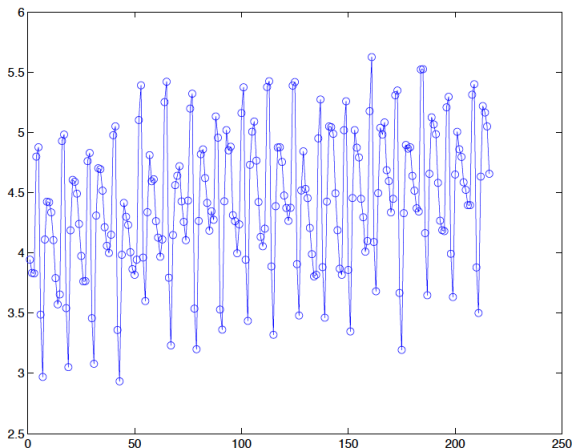
Time-series transformations - Rolling mean



Compact time-series description

$$X(t) = \text{trend}(t) + \text{seasonal}(t) + \text{residual}(t)$$

- **Example:** Suppose we have the following time-series:

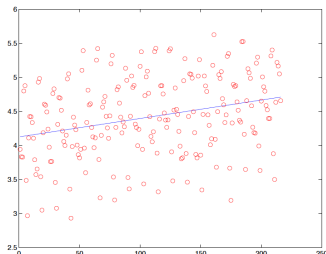


Compact time-series description

- In terms of equations, (roughly) we perform two types of regression, and obtain residuals

$$X(t) = \text{polynomial}(t) + \sum_i (\beta_i \cos(\lambda_i t) + \gamma_i \sin(\lambda_i t)) + R_t$$

- First, we perform regression: $X(t) = \alpha_0 + \alpha_1 t + E(t)$:



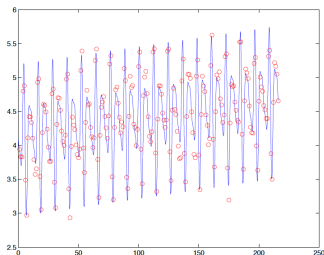
Compact time-series description

- In terms of equations, (roughly) we perform two types of regression, and obtain residuals

$$X(t) = \text{polynomial}(t) + \sum_i (\beta_i \cos(\lambda_i t) + \gamma_i \sin(\lambda_i t)) + R_t$$

- Second regression

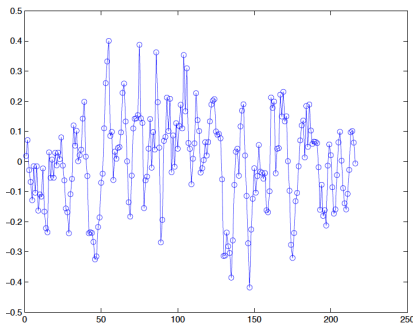
$$X(t) - (\alpha_0 + \alpha_1 t) = \sum_i (\beta_i \cos(\lambda_i t) + \gamma_i \sin(\lambda_i t)) + E(t):$$



Compact time-series description

- Left with the residuals:

$$R(t) = X(t) - (\alpha_0 + \alpha_i t) - \sum_i (\beta_i \cos(\lambda_i t) + \gamma_i \sin(\lambda_i t))$$



- **Not done yet!** Study the residuals. E.g., can we forecast them?

Compact time-series description

- In Python things have been made easy for us.

```
1 def _decom(dataset):
2     decomposition = statsmodels.tsa.seasonal.
3     seasonal_decompose(dataset)
4     plt.plot(decomposition.trend, label='Trend',
5             color = 'b')
6     plt.plot(decomposition.seasonal, label='
7     Seasonality', color = 'b')
8     plt.plot(decomposition.resid, label='
9     Residuals', color = 'b')
```

Today's outline

① Part I: Fundamentals

- Stochastic processes, strict vs weak stationarity, correlograms, partial correlograms, periodograms, autoregressive models (AR), moving average models (MA), transformations, compact time-series description.

② Part II: Forecasting methods

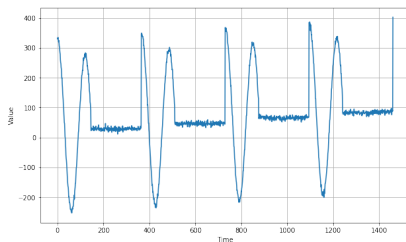
- Evaluation protocols, basic metrics, null models, spectral methods, ARMA, ARIMA, SARIMA, VAR, similarity search, deep learning (DNNs, RNNs)

Train and test – Respect order!

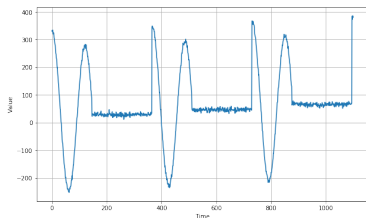


```
1 split_time = 1000
2 x_train = series[:split_time]
3 x_test = series[split_time:]
```

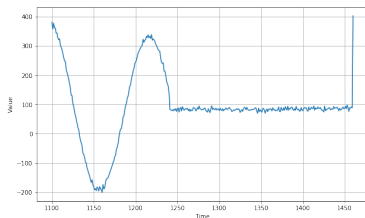
Train and test – Respect order!



Dataset



Training set



Test set

Evaluation protocols

There exist **two main evaluation protocols** for a forecasting method.

① k -fold cross validation

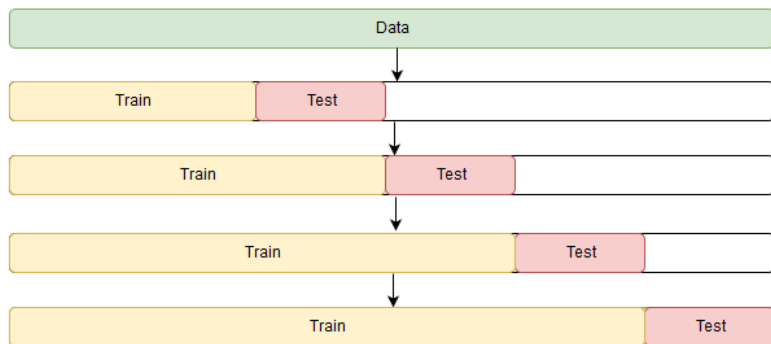
Remarks

- Cannot just randomly partition the time-series into folds as in other ML tasks (**respect the order!**)

② Walk forward evaluation protocol

- Both protocols can be easily explained visually.

k-fold cross validation



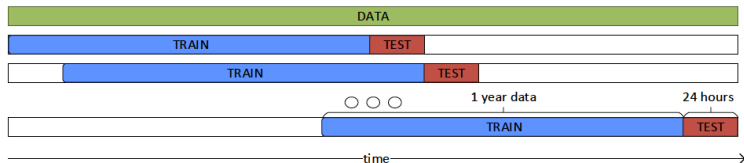
- Test set size remains same, but training size increases.

```
1 from sklearn.model_selection import  
   TimeSeriesSplit
```

k-fold cross validation – sklearn

```
1 X = np.array([[1, 2], [3, 4], [1, 2], [3, 4],  
2             [1, 2], [3, 4]])  
3 y = np.array([1, 2, 3, 4, 5, 6])  
4 tscv = TimeSeriesSplit(max_train_size=None,  
5                       n_splits=5)  
6 for train_index, test_index in tscv.split(X):  
7     ...     print("TRAIN:", train_index, "TEST:",  
8             test_index)  
9     ...     X_train, X_test = X[train_index], X[  
10            test_index]  
11    ...     y_train, y_test = y[train_index], y[  
12            test_index]  
13 >> TRAIN: [0] TEST: [1]  
14 >> TRAIN: [0] TEST: [1]  
15 >> TRAIN: [0 1] TEST: [2]  
16 >> TRAIN: [0 1 2] TEST: [3]  
17 >> TRAIN: [0 1 2 3] TEST: [4]  
18 >> TRAIN: [0 1 2 3 4] TEST: [5]
```


Walk-forward evaluation protocol



- Sliding window of same size.
- Straight-forward implementation.
- In reality, when training is expensive we may sample as many as possible such windows, with a bias towards more recent windows.

Basic metrics – g groundtruth, f forecast

```
1 import numpy as np
2 errors = forecast - truth
```

- Mean squared error (mse):

$$mse(f, g) = \frac{100\%}{n} \sum_{t=1}^n (f_t - g_t)^2.$$

```
1 mse = np.square(errors).mean()
2
```

- Root mean squared error (rmse):

$$rmse(f, g) = \sqrt{mse(f, g)}.$$

```
1 rmse = np.sqrt(mse)
2
```

Basic metrics

- Mean absolute error (mae):

$$mae(f, g) = \frac{100\%}{n} \sum_{t=1}^n |(f_t - g_t)|.$$

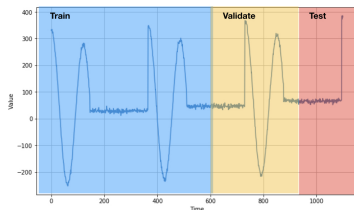
```
1 mae = np.abs(errors).mean()  
2
```

- Mean absolute percentage error (mape):

$$mape(f, g) = \frac{100\%}{n} \sum_{t=1}^n \left| \frac{f_t - g_t}{g_t} \right|.$$

```
1 mape = np.abs(errors/x_test).mean()  
2
```

Train, validate and test



- Alternatively, we can split the data into
 - ① Training set
 - ② Validation set in order to evaluate hyperparameters, useful for minimizing overfitting
 - ③ Test set

```
1 x_train = series[:split1]
2 x_validate = series[split1:split2]
3 x_test = series[split2:]
```

Naive forecasting

- Suppose we wish to predict $x(t + 1)$ given past values (and possibly other time-series).
- **Naive forecasting:**

$$x(t + 1) \leftarrow x(t).$$

- **Example:** Suppose we wish to predict Torino's temperature next hour.
Prediction: The temperature in the next hour is going to be the same as now.
- Competitive baseline!

More null models

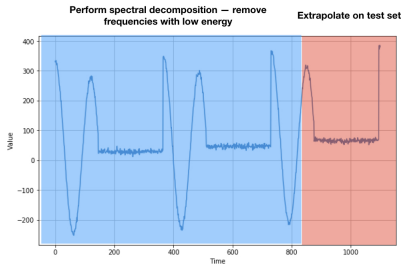
- Constant values, e.g., global average or a frequently occurring value in the dataset.
- Average/median of past most recent w values where w can be small or spanning the whole dataset

$$x(t+1) \leftarrow \frac{\sum_{j=1}^w x(t-j)}{w}.$$

- **Remark:** $w = 1$, i.e., naive forecasting, tends to be the most competitive baseline.
- Exponential smoothing (more importance to most recent measurements)
- Regression and extrapolation
- ...

Spectral forecasting

- ① We pretend the training dataset is periodic.
- ② We perform a spectral decomposition of the signal in sinusoidal functions.
- ③ We retain the frequencies that hold most but **not all** of the energy
- ④ **Rule of thumb:** keeping 80% of the energy is a good choice, but more values should be tried out.
- ⑤ Extrapolate the periodic signal to the test set.



Spectral forecasting – Training

```
1 def spectral_prediction(timeseries, n_predict,
2   perc = 0.8):
3     n = timeseries.size
4     x_freqdom = fft(timeseries)
5     f = np.fft.fftfreq(n)
6     indexes = list(range(n))
7     harm_squares = np.square(np.abs(x_freqdom))
8     sum_squares = harm_squares.sum()
9     sortd = np.argsort(-harm_squares)
10    cum_sum = 0
11    i = 0
12    n_harm = 0 # number of harmonics in model,
13    keep perc% of energy
14    while cum_sum < sum_squares*perc:
15        cum_sum += harm_squares[sortd[i]]
16        i+=1
17        n_harm+=1
```


Spectral forecasting – Training

```
1     my_signal = np.zeros(t.size)
2     for i in indexes:
3         # This is the amplitude of the
4         extrapolated frequency
5         ampli = np.absolute(x_freqdom[i]) / n
6         # This is the phase of the extrapolated
7         frequency
8         phase = np.angle(x_freqdom[i])
9         # we add this component to the
10        extrapolated signal
11        my_signal += ampli * np.cos(2 * np.pi *
12        f[i] * t + phase)
13
14    return my_signal
```

ARMA

- An ARMA(p, q) process $\{X_t\}$ is a stationary process that satisfies

$$X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p} = W_t + \theta_1 W_{t-1} + \dots + W_{t-q},$$

where $W_t \sim WN(0, \sigma^2)$.

- **Observations**
 - ① ARMA($p, 0$) is same as AR(p).
 - ② ARMA($0, q$) is same as MA(q).

ARMA

- ARMA processes are very important because of the following (informally) stated theorem:

Theorem (Informal statement)

For any stationary process $\{Y_t\}$ with autocovariance γ , and any $k > 0$ there is an ARMA process $\{X_t\}$ that fits $\{Y_t\}$ well.

- This is one **big** reason why we try to transform a time-series into a stationary process.
- Many null models are special cases of this model.
- Interpretable.
- For certain types of data it is performing well but for volatile, non-stationary time-series it is not.
- **Always** fundamental, non-trivial baseline.

ARMA - Lowest terms

Consider a process $\{X_t\}$ such that $X_t = W_t$ where $W_t \sim WN(0, \sigma^2)$.

$$X_t - 3X_{t-1} + 4X_{t-2} = W_t - 3W_{t-1} + 4W_{t-2}$$

- This looks like an ARMA(2,2) process! But in reality it is not, since it is not.
- Something called *characteristic polynomials* are not in lowest terms, they share common factors. Actually, here both the LHS and the RHS have the same characteristic polynomial $\phi(B) = 1 - 3B + 4B^2$.
- Lowest terms means that the following fraction cannot be not be simplified further

$$\frac{\phi_{RHS}(B)}{\phi_{LHS}(B)}$$

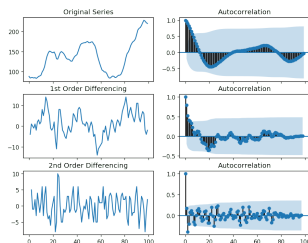
ARIMA

- ARIMA stands for “Auto Regressive Integrated Moving Average”
- In addition to parameters p, q we have another parameter d
- d is the number of differencing required to make the time series stationary

E.g., when $d = 0$ we have a standard ARMA model.

Remark: differencing does not always succeed, and overdifferencing may produce “bad” time-series.

ARIMA



- ARIMA extends to SARIMA by including a seasonal component.
- This involves also p, d, q parameters but also the seasonality parameter s .

```
1 model=sm.tsa.statespace.SARIMAX( data['PSV -  
   Italian gas\n[ /MWh]'], order=(1,1,1),  
2  
   seasonal_order=(1,1,1,1),  
   enforce_invertibility=False).fit(dispatch=-1)  
3 fitted1 = model.fittedvalues
```

SARIMA - Parameter search

- ACF, PACF plots but also exhaustive search

```
1 def parameter_search():
2     p = range(1, 11), q = range(1, 11)
3     d=range(1,3)
4     Ps = [1,2,3], D= range(1,3), Qs = [1,2,3]
5     s = 7
6     parameters = product(p,d,q, Ps, D, Qs)
7     parameters_list = list(parameters)
8     results = []
9     best_aic = float("inf")
10    for param in tqdm(parameters_list):
11        model=SARIMAX(data, order=(param[0],
12        param[1], param[2]), seasonal_order=(param
13        [3], param[4], param[5], s)).fit(dispatch=-1)
14        aic = model.aic
15        if aic < best_aic:
16            print('Best model so far :', param)
```

SARIMA(X) – Summary

- Summary
 - ① Very important family of forecasting methods
 - ② Well developed packages in Python and R
 - ③ Lots of theory behind them
 - ④ Under stationarity, they come with lots of nice properties (e.g., confidence intervals)
 - ⑤ However stationarity is not always there...
 - ⑥ Therefore, it is always a baseline.
 - ⑦ SARIMAX also supports exogenous variables, hence the X.

VAR

- For multivariate series, these concepts have natural generalizations.
- For example, the next equation defines a vector autoregressive process of order 1.

$$\begin{bmatrix} y_{1,t} \\ y_{2,t} \end{bmatrix} = \begin{bmatrix} a_{1,1} & a_{1,2} \\ a_{2,1} & a_{2,2} \end{bmatrix} \cdot \begin{bmatrix} y_{1,t-1} \\ y_{2,t-1} \end{bmatrix} + \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \quad (1)$$

- Python *statsmodels* package has VAR ready for us.

```
1 from statsmodels.tsa.vector_ar.var_model import
   VAR
2 model = VAR(endog=train)
3 myfit = model.fit()
4 prediction = myfit.forecast(y=data[-1,:], steps
   =1)
```

Similarity search

- When are two time series $x(t), y(t), t = 1, \dots, n$ similar?

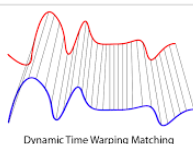
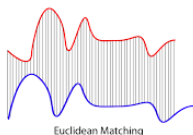
How do we quantify their similarity?

- There exist **two** major families of distances:

① Euclidean and ℓ_p norms

- Euclidean distance $\sum_{t=1}^n (x(t) - y(t))^2$ (ℓ_2 distance)
- Manhattan distance $\sum_{t=1}^n |x(t) - y(t)|$ (ℓ_1 distance)
- Dot product $\langle x, y \rangle$
- ...

② Time warping (DTW) and variations



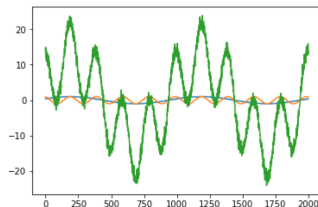
Similarity search

- In Python, again things are pretty straight-forward.

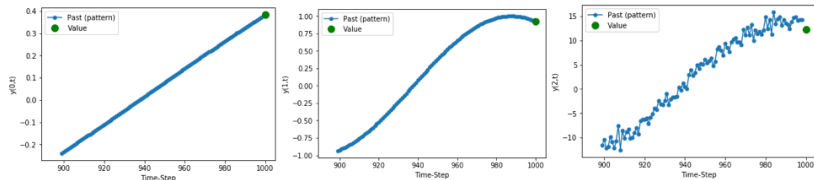
```
1 def euclid(points):
2     return np.linalg.norm(points[0]-points[1])
3
4 def l1_distance(points):
5     return np.linalg.norm(points[0]-points[1],
6                             ord=1)
7
8 from dtw import dtw
9 d, cost_matrix, acc_cost_matrix, path = dtw(x, y
10      , dist=l1_distance)
```

Similarity search - Toy dataset

```
1 def demo1():
2     x = np.linspace(-np.pi, np.pi, 2001)
3     s1= np.sin(2*x+np.pi/8)
4     s2= np.cos(10*x+np.pi/8)
5     target = 12*s1+10*s2+np.random.normal(loc
6         =0.0, scale=1, size=len(x))
7     dataset = pd.DataFrame({'f1': s1, 'f2': s2,
8         'target':target})
9     df = dataset
10    data = df.values.T
11    return data
```



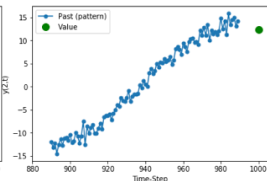
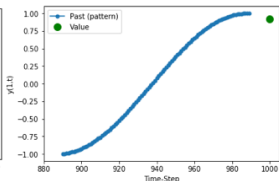
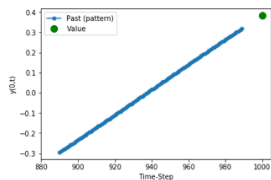
Similarity search - Windowing around timestamp



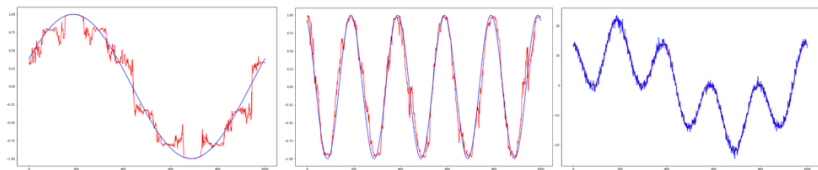
- We simultaneously predict all time series, not just the target time series
- We look into the w previous values and create a vector with $dims \times w$ coordinates.
- We search for the $k \geq 1$ most similar such vector in the training dataset (k small).
- We output the mean/median of the top- k nearest neighbors.

Similarity search

- Sometimes, we don't have the immediate w past values available.
- Then, we just use the w most recent past values that we have available for each time-series.

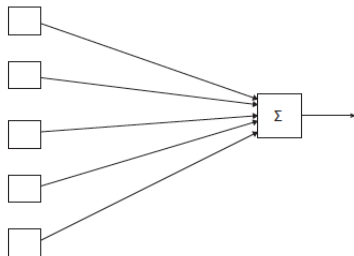
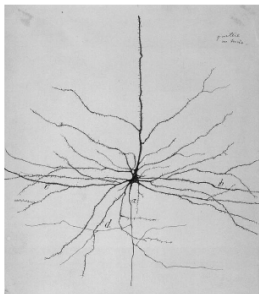


Similarity search – Prediction results



- In this toy example the prediction results are good overall.
- Similarity search is a useful framework that is able to either provide decent predictions or show why other methods fail to product good outputs
 - i.e., lots of variance among similar windowed patterns.
- Finding the “right” notion of distance is an important component, as well as the right window length.
- For high-dimensional search, LSH can be used.

Neurons and perceptrons



A human brain neuron and a perceptron

Tensorflow – Open source for deep learning



Tensorflow

- Tensorflow is an end-to-end open source machine learning platform...
- but is also a symbolic math library that can be used even to run a favorite optimization algorithm on your problem.

```
1 import tensorflow as tf
2 x = tf.constant("Hello world")
3 sess = tf.Session()
4 print(sess.run(x))
```

Deep learning

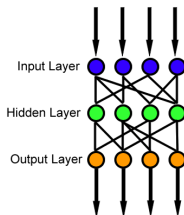
- We have a bunch of inputs associated with an output.
- Can we see the right output for a new input, i.e., an input we have not seen?
- **Deep learning** is a leading approach to data science, that uses multiple layers to automatically extract high-quality, higher-level features from the raw inputs.
 - deep neural networks, recurrent neural networks, convolutional neural networks etc.
- **Why does it work?** Increasing amount of research to understand why does it work, even when neurons use simple activation functions (piecewise flat.)

Deep feedforward neural networks (DNNs)

Deep feedforward networks feedforward, also often called neural networks, or multilayer perceptrons (MLPs), are the quintessential deep learning models.

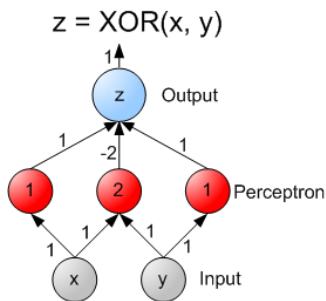
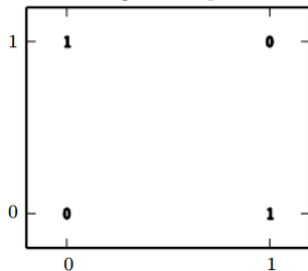
- The goal of a feedforward network is to approximate some function f^* by composing layers, e.g.,

$$\tilde{f}(x) = f^{(3)}(f^{(2)}(f^{(1)}(x))).$$



Deep feedforward neural networks (DNNs) – Example

- Classic example: DNNs can learn the XOR function.



DNNs for time-series – Step 1

1. Prepare time series data for DNN by creating a windowed dataset:
 - For each timestamp, let $x(t)$ be the value/label.
 - The previous w values could be seen as the input features.
 - Lots of space for feature engineering (Fourier coeffs, wavelet coeffs, min/max/median etc.)

```
1 dataset = tf.data.Dataset.range(20)
2 dataset = dataset.window(5, shift=1,
3   drop_remainder=True)
4 dataset = dataset.flat_map(lambda window: window
5   .batch(5))
6 dataset = dataset.map(lambda window: (window
7  [:-1], window[-1:]))
8 for x,y in dataset:
9     print(x.numpy(), y.numpy())
```

DNNs for time-series – Step 1

```
1 [0 1 2 3] [4]
2 [1 2 3 4] [5]
3 [2 3 4 5] [6]
4 [3 4 5 6] [7]
5 [4 5 6 7] [8]
6 [5 6 7 8] [9]
7 [6 7 8 9] [10]
8 [ 7 8 9 10] [11]
9 [ 8 9 10 11] [12]
10 [ 9 10 11 12] [13]
11 [10 11 12 13] [14]
12 [11 12 13 14] [15]
13 [12 13 14 15] [16]
14 [13 14 15 16] [17]
15 [14 15 16 17] [18]
16 [15 16 17 18] [19]
```

DNNs for time-series – Step 1b (batching)

```
1 dataset = dataset.shuffle()
2 dataset = dataset.batch(2).prefetch(1)
3 x = [[5 6 7 8] [1 2 3 4]]
4 y = [[9] [5]]
5 x = [[ 9 10 11 12] [12 13 14 15]]
6 y = [[13] [16]]
7 x = [[11 12 13 14] [13 14 15 16]]
8 y = [[15] [17]]
9 x = [[10 11 12 13] [ 4 5 6 7]]
10 y = [[14] [ 8]]
11 x = [[ 7 8 9 10] [14 15 16 17]]
12 y = [[11] [18]]
13 x = [[3 4 5 6] [0 1 2 3]]
14 y = [[7] [4]]
15 x = [[15 16 17 18] [ 6 7 8 9]]
16 y = [[19] [10]]
17 x = [[ 8 9 10 11] [ 2 3 4 5]]
18 y = [[12] [ 6]]
19
```

DNNs for time-series – Step 2

2. Split the data into training and test sets.

```
1 time = np.arange(3 * 365, dtype="float32")
2 split_time = 2*365+31+28+30
3 time_train = time[:split_time]
4 x_train = series[:split_time]
5 time_test = time[split_time:]
6 x_test = series[split_time:]
```


DNNs for time-series – Steps 3, 4

3. Setup a DNN architecture (here single layer)

4. Train it

4b Inspect the layer weights

```
1 dataset = windowed_dataset(x_train, window_size,
    batch_size)
2 layer = tf.keras.layers.Dense(1, input_shape=[
    window_size])
3 model = tf.keras.models.Sequential([layer])
4 model.compile(loss="mse", optimizer=tf.keras.
    optimizers.SGD(lr=1e-6, momentum=0.9))
5 model.fit(dataset, epochs=100, verbose=0)
6 print("Layer weights {}".format(layer.
    get_weights()))
```

DNNs for time-series – Steps 3, 4

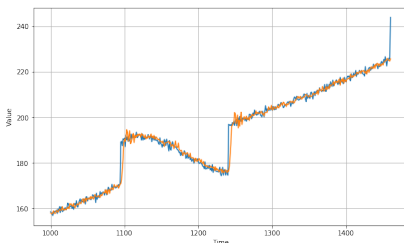
```
1 Layer weights [array([[ -0.03136637],
2     [ -0.04519343],
3     [ 0.05885418],     [ 0.05017925],
4     [-0.01355846],     [-0.07844295],
5     [ 0.04997651],     [ 0.03857101],
6     [-0.01163514],     [ 0.02050608],
7     [-0.01307715],     [-0.05444159],
8     [ 0.01337368],     [ 0.05073489],
9     [ 0.01348611],     [-0.02583448],
10    [ 0.08094374],     [ 0.23290157],
11    [ 0.20300445],     [ 0.46333405]],
    dtype=float32), array([0.01811779], dtype=
    float32)]
```

- These weights are what we need to forecast a value given a new input of 20 values.

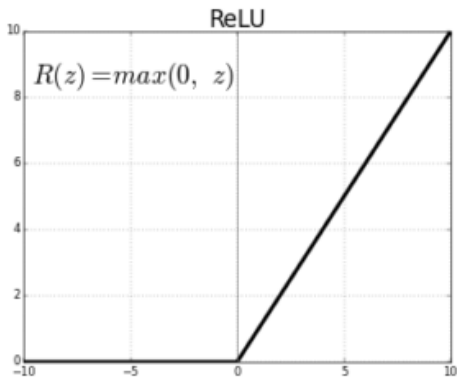
DNNs for time-series – Step 5

5. Forecast

```
1 forecast = []
2 for time in range(len(series) - window_size):
3     forecast.append(model.predict(series[time:time
4         + window_size][np.newaxis]))
5 forecast = forecast[split_time-window_size:]
6 plt.plot(time_test, forecasts)
7 tf.keras.metrics.mean_absolute_error(x_valid,
    results).numpy()
```



Quick comment - Activation function ReLU



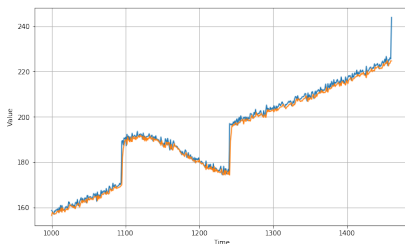
- ReLU is an important activation function, especially because it is simple, yet sparse-friendly.

DNNs for time-series – Adding layers

- Our example uses one layer. We can easily add layers as follows.

```
1 model = tf.keras.models.Sequential([tf.keras.layers.Dense(10, input_shape=[window_size],  
activation="relu"), tf.keras.layers.Dense(10,  
activation="relu"), tf.keras.layers.Dense(1)  
])
```

2



DNNs for time-series – Did the extra layers help?

- In addition to the visualization, KERAS provides an easy way to compute the various metrics.

```
1 tf.keras.metrics.mean_absolute_error(x_test ,
   results_shallow).numpy()
2 >> 1.412784
3 tf.keras.metrics.mean_absolute_error(x_test ,
   results_deep).numpy()
4 >>1.3640015
```

- It is worth outlining that this is not always the case.
- The number of layers is one of the hyperparameters to be optimized.

Remarks on DNNs

- For multivariate time series, we change the input shape from a vector to a matrix (batch size \times number of time-series/dimensions)
- Lots of other parameters need to be fine-tuned, including the weight initializer, use or not of batch normalization, dropout rates, learning rates of SGD etc.
- Feature selection earlier is also useful for DNNs.
- Not suited for sequential data, but for supervised framework.

- Next: Deep RNNs

Last important remark about DNNs

- The approach we described for analyzing time-series can be extended in two different directions.
- ① Create more features from the window. For instance add dimensions for the maximum value observed.

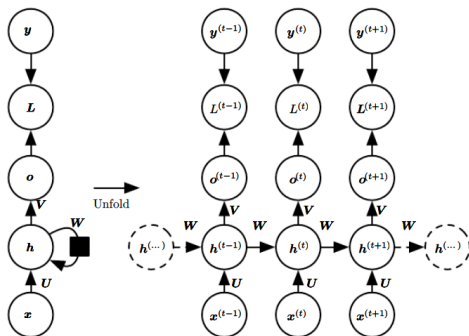
```
1 from tsfresh import extract_features
```

tsfresh is a library that allows for fast extraction of hundreds of features from each window.

- ② We can use any other supervised ML method including SVM regressors, regression random forests, Gaussian processes etc.

Deep recursive neural networks (RNNs)

- RNNs are neural networks for sequential data.



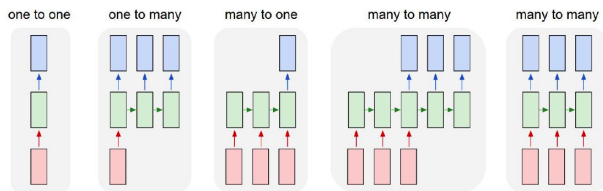
Source: Deep Learning book

$$h_t = f_W(h_{t-1}, x_t).$$

- Same function and same set of parameters W are used at each time step.

Deep recursive neural networks (RNNs)

- They have a distributed hidden state that allows them to store a information about the past.
- RNNs offer a lot of flexibility that makes them valuable in diverse applications



Source: Andrej Karpathy

- An RNN can be thought of as multiple copies of the same network, each passing a message to a successor.

Deep recursive neural networks (RNNs) – TensorFlow

- In contrast to DNNs, the input shape is

[batch size, #timesteps, #dimensions]

```
1 model = tf.keras.models.Sequential([
2     tf.keras.layers.SimpleRNN(40, return_sequences
3         =True),
4     tf.keras.layers.SimpleRNN(40),
5     tf.keras.layers.Dense(1)
6 ])
7 optimizer = tf.keras.optimizers.SGD(lr=5e-5,
8     momentum=0.9)
9 model.compile(loss=tf.keras.losses.Huber(),
10     optimizer=optimizer,
11     metrics=["mae"])
12 model.fit(dataset, epochs=100)
```

Deep recursive neural networks (RNNs) – Tensorflow

- Prediction is done in the same way, and evaluation too, e.g., using

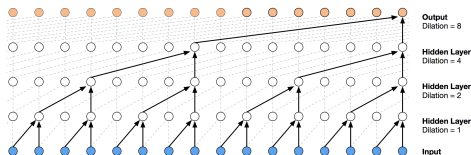
```
1 tf.keras.metrics.mean_absolute_error(x_test,  
   forecast_results).numpy()
```

- The key method is MODEL.PREDICT again:

```
1 forecast=[]  
2 for time in range(len(series) - window_size):  
3     forecast.append(model.predict(series[time:time  
   + window_size][np.newaxis]))  
4  
5 forecast = forecast[split_time-window_size:]
```

Remark: LSTMs and dilated convolutions

- RNNs are supposed to remember from past, but in practice they forget easily.
- LSTMs are units that are able to capture better long term dependencies.
- State-of-the-art time-series prediction methods combine convolutions, typically used in image processing, with recurrent neural networks.



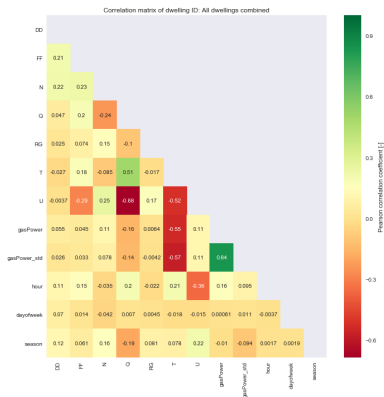
Dilated convolutions

Source: original paper, *Temporal convolutional networks*
by Bai et al. Arxiv:1803.01271



Tip of iceberg – Lots more to talk about!

- Today I have scratched the tip of the iceberg.
- E.g., in multivariate analysis how do we select the most relevant time-series to the target series we wish to forecast?



Thank you! Questions?

web page: <http://tsourakakis.com>

email: ctsourak@bu.edu

