

Time Series Analysis with Python

Aileen Nielsen

July, 12, 2016

aileen.a.nielsen@gmail.com

Installation instructions

- Please install Conda per 'quick install' instructions:
<http://conda.pydata.org/docs/install/quick.html>
- Make sure you have the following packages installed:
 - pandas
 - numpy
 - Statsmodels
 - scikit-learn
 - scipy
- These would be good to have but are not essential:
 - pytz
 - hmm

Dealing with time

Pandas functionality

- Generate sequences of fixed-frequency dates and time spans
- Conform or convert time series to a particular frequency
- Compute 'relative' dates based on various non-standard time increments (e.g. 5 business days before the last day of the year) or 'roll' dates backward and forward

Date/Time components

Time/Date Components

There are several time/date properties that one can access from `Timestamp` or a collection of timestamps like a `DateTimeIndex`.

Property	Description
<code>year</code>	The year of the datetime
<code>month</code>	The month of the datetime
<code>day</code>	The days of the datetime
<code>hour</code>	The hour of the datetime
<code>minute</code>	The minutes of the datetime
<code>second</code>	The seconds of the datetime
<code>microsecond</code>	The microseconds of the datetime
<code>nanosecond</code>	The nanoseconds of the datetime
<code>date</code>	Returns <code>datetime.date</code>
<code>time</code>	Returns <code>datetime.time</code>
<code>dayofyear</code>	The ordinal day of year
<code>weekofyear</code>	The week ordinal of the year
<code>week</code>	The week ordinal of the year
<code>dayofweek</code>	The numer of the day of the week with Monday=0, Sunday=6
<code>weekday</code>	The number of the day of the week with Monday=0, Sunday=6
<code>weekday_name</code>	The name of the day in a week (ex: Friday)
<code>quarter</code>	Quarter of the date: Jan=Mar = 1, Apr-Jun = 2, etc.
<code>days_in_month</code>	The number of days in the month of the datetime
<code>is_month_start</code>	Logical indicating if first day of month (defined by frequency)
<code>is_month_end</code>	Logical indicating if last day of month (defined by frequency)
<code>is_quarter_start</code>	Logical indicating if first day of quarter (defined by frequency)
<code>is_quarter_end</code>	Logical indicating if last day of quarter (defined by frequency)
<code>is_year_start</code>	Logical indicating if first day of year (defined by frequency)
<code>is_year_end</code>	Logical indicating if last day of year (defined by frequency)

Offset Aliases

A number of string aliases are given to useful common time series frequencies. We will refer to these aliases as *offset aliases* (referred to as *time rules* prior to v0.8.0).

Alias	Description
B	business day frequency
C	custom business day frequency (experimental)
D	calendar day frequency
W	weekly frequency
M	month end frequency
BM	business month end frequency
CBM	custom business month end frequency
MS	month start frequency
BMS	business month start frequency
CBMS	custom business month start frequency
Q	quarter end frequency
BQ	business quarter end frequency
QS	quarter start frequency
BQS	business quarter start frequency
A	year end frequency
BA	business year end frequency
AS	year start frequency
BAS	business year start frequency
BH	business hour frequency
H	hourly frequency
T, min	minutely frequency
S	secondly frequency
L, ms	milliseconds
U, us	microseconds
N	nanoseconds

DateOffset objects

In the preceding examples, we created `DatetimeIndex` objects at various frequencies by passing in frequency strings like 'M', 'W', and 'BM' to the `freq` keyword. Under the hood, these frequency strings are being translated into an instance of pandas `DateOffset`, which represents a regular frequency increment. Specific offset logic like "month", "business day", or "one hour" is represented in its various subclasses.

Class name	Description
<code>DateOffset</code>	Generic offset class, defaults to 1 calendar day
<code>BDay</code>	business day (weekday)
<code>CDay</code>	custom business day (experimental)
<code>Week</code>	one week, optionally anchored on a day of the week
<code>WeekOfMonth</code>	the x-th day of the y-th week of each month
<code>LastWeekOfMonth</code>	the x-th day of the last week of each month
<code>MonthEnd</code>	calendar month end
<code>MonthBegin</code>	calendar month begin
<code>BMonthEnd</code>	business month end
<code>BMonthBegin</code>	business month begin
<code>CBMonthEnd</code>	custom business month end
<code>CBMonthBegin</code>	custom business month begin
<code>QuarterEnd</code>	calendar quarter end
<code>QuarterBegin</code>	calendar quarter begin
<code>BQuarterEnd</code>	business quarter end
<code>BQuarterBegin</code>	business quarter begin
<code>FY5253Quarter</code>	retail (aka 52-53 week) quarter
<code>YearEnd</code>	calendar year end
<code>YearBegin</code>	calendar year begin
<code>BYearEnd</code>	business year end
<code>BYearBegin</code>	business year begin
<code>FY5253</code>	retail (aka 52-53 week) year
<code>BusinessHour</code>	business hour
<code>CustomBusinessHour</code>	custom business hour
<code>Hour</code>	one hour
<code>Minute</code>	one minute
<code>Second</code>	one second
<code>Milli</code>	one millisecond
<code>Micro</code>	one microsecond
<code>Nano</code>	one nanosecond

Anchored offsets

Anchored Offsets

For some frequencies you can specify an anchoring suffix:

Alias	Description
W-SUN	weekly frequency (sundays). Same as 'W'
W-MON	weekly frequency (mondays)
W-TUE	weekly frequency (tuesdays)
W-WED	weekly frequency (wednesdays)
W-THU	weekly frequency (thursdays)
W-FRI	weekly frequency (fridays)
W-SAT	weekly frequency (saturdays)
(B)Q(S)-DEC	quarterly frequency, year ends in December. Same as 'Q'
(B)Q(S)-JAN	quarterly frequency, year ends in January
(B)Q(S)-FEB	quarterly frequency, year ends in February
(B)Q(S)-MAR	quarterly frequency, year ends in March

Indices

DatetimeIndex

One of the main uses for `DatetimeIndex` is as an index for pandas objects. The `DatetimeIndex` class contains many timeseries related optimizations:

- A large range of dates for various offsets are pre-computed and cached under the hood in order to make generating subsequent date ranges very fast (just have to grab a slice)
- Fast shifting using the `shift` and `tshift` method on pandas objects
- Unioning of overlapping `DatetimeIndex` objects with the same frequency is very fast (important for fast data alignment)
- Quick access to date fields via properties such as `year`, `month`, etc.
- Regularization functions like `snap` and very fast `asof` logic

Pandas time-relates data types

Overview

Following table shows the type of time-related classes pandas can handle and how to create them.

Class	Remarks	How to create
Timestamp	Represents a single time stamp	<code>to_datetime</code> , <code>Timestamp</code>
DatetimeIndex	Index of <code>Timestamp</code>	<code>to_datetime</code> , <code>date_range</code> , <code>DatetimeIndex</code>
Period	Represents a single time span	<code>Period</code>
PeriodIndex	Index of <code>Period</code>	<code>period_range</code> , <code>PeriodIndex</code>

Holidays / Holiday Calendars

Holidays and calendars provide a simple way to define holiday rules to be used with `CustomBusinessDay` or in other analysis that requires a predefined set of holidays. The `AbstractHolidayCalendar` class provides all the necessary methods to return a list of holidays and only `rules` need to be defined in a specific holiday calendar class. Further, `start_date` and `end_date` class attributes determine over what date range holidays are generated. These should be overwritten on the `AbstractHolidayCalendar` class to have the range apply to all calendar subclasses. `USFederalHolidayCalendar` is the only calendar that exists and primarily serves as an example for developing other calendars.

For holidays that occur on fixed dates (e.g., US Memorial Day or July 4th) an observance rule determines when that holiday is observed if it falls on a weekend or some other non-observed day. Defined observance rules are:

Rule	Description
<code>nearest_workday</code>	move Saturday to Friday and Sunday to Monday
<code>sunday_to_monday</code>	move Sunday to following Monday
<code>next_monday_or_tuesday</code>	move Saturday to Monday and Sunday/Monday to Tuesday
<code>previous_friday</code>	move Saturday and Sunday to previous Friday"
<code>next_monday</code>	move Saturday and Sunday to following Monday

An example of how holidays and holiday calendars are defined:

Thorns in your side

pandas.DataFrame.asfreq

`DataFrame.asfreq(freq, method=None, how=None, normalize=False)`

Convert all TimeSeries inside to specified frequency using DateOffset objects. Optionally provide fill method to pad/backfill missing values.

Parameters:

freq : *DateOffset object, or string*

method : {'backfill', 'bfill', 'pad', 'ffill', None}

Method to use for filling holes in reindexed Series pad / ffill: propagate last valid observation forward to next valid backfill / bfill: use NEXT valid observation to fill method

how : {'start', 'end'}, default end

For PeriodIndex only, see PeriodIndex.asfreq

normalize : bool, default False

Whether to reset output index to midnight

Returns:

converted : *type of caller*

pandas.DataFrame.resample

`DataFrame.resample(rule, how=None, axis=0, fill_method=None, closed=None, label=None, convention='start', kind=None, loffset=None, limit=None, base=0)`

Convenience method for frequency conversion and resampling of regular time-series data.

Parameters:

rule : *string*

the offset string or object representing target conversion

axis : *int, optional, default 0*

closed : {'right', 'left'}

Which side of bin interval is closed

label : {'right', 'left'}

Which bin edge label to label bucket with

convention : {'start', 'end', 's', 'e'}

loffset : *timedelta*

Adjust the resampled time labels

base : *int, default 0*

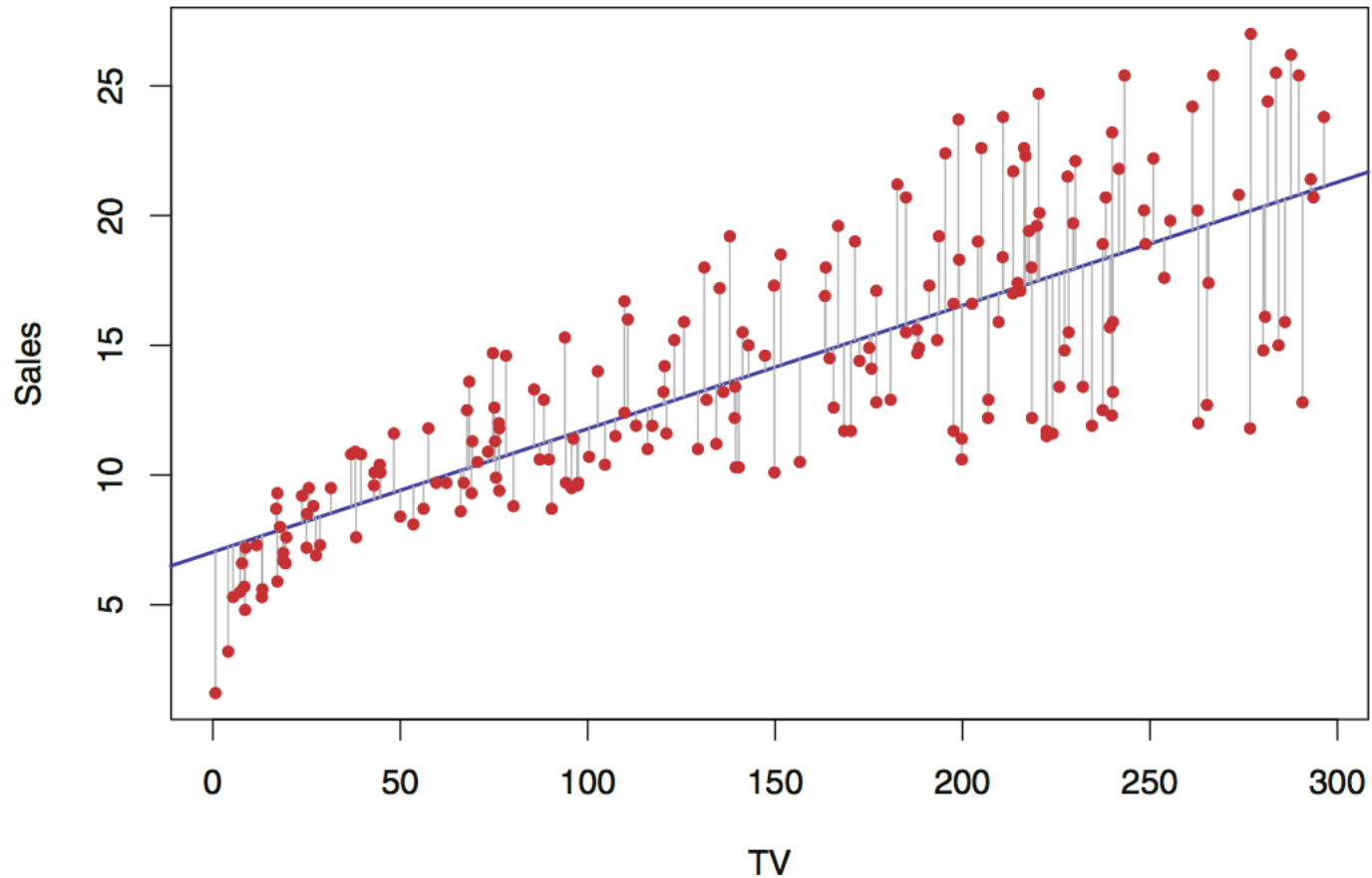
For frequencies that evenly subdivide 1 day, the "origin" of the aggregated intervals.

For example, for '5min' frequency, base could range from 0 through 4. Defaults to 0

Examples

Linear regression

Linear regression intuition

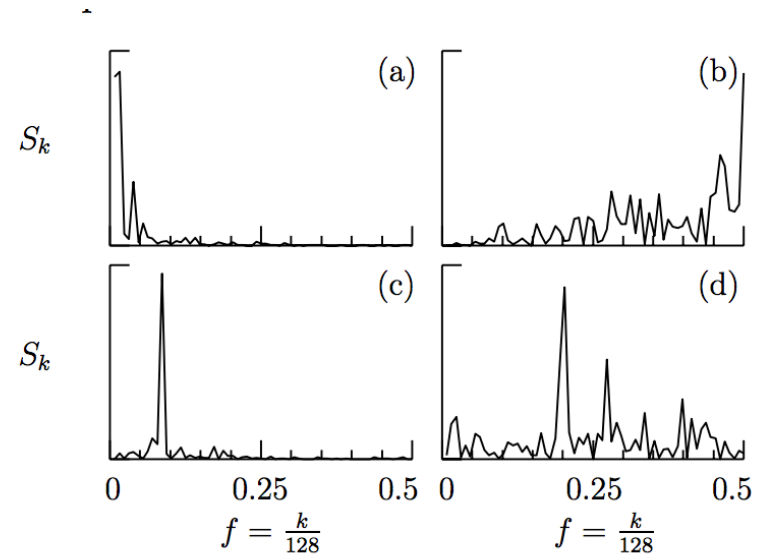
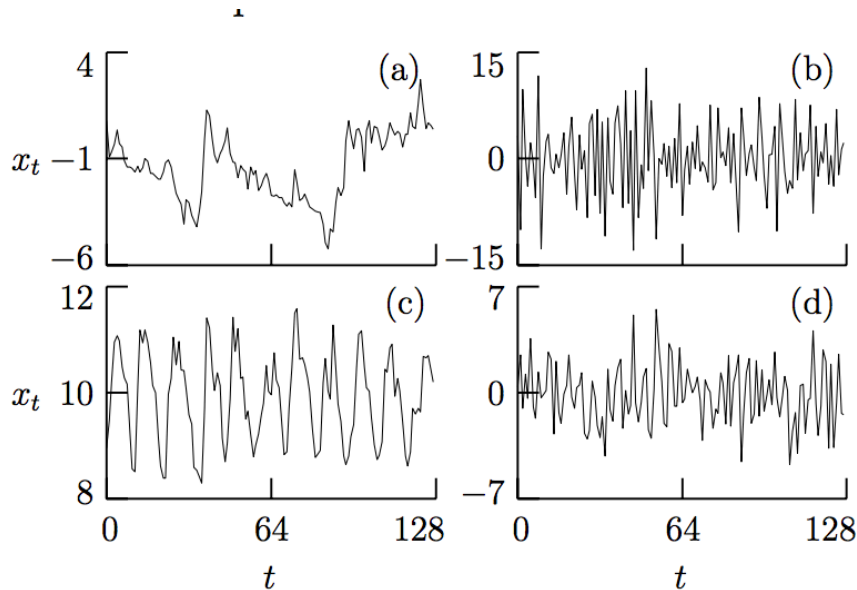


Spectral Analysis

Intuition

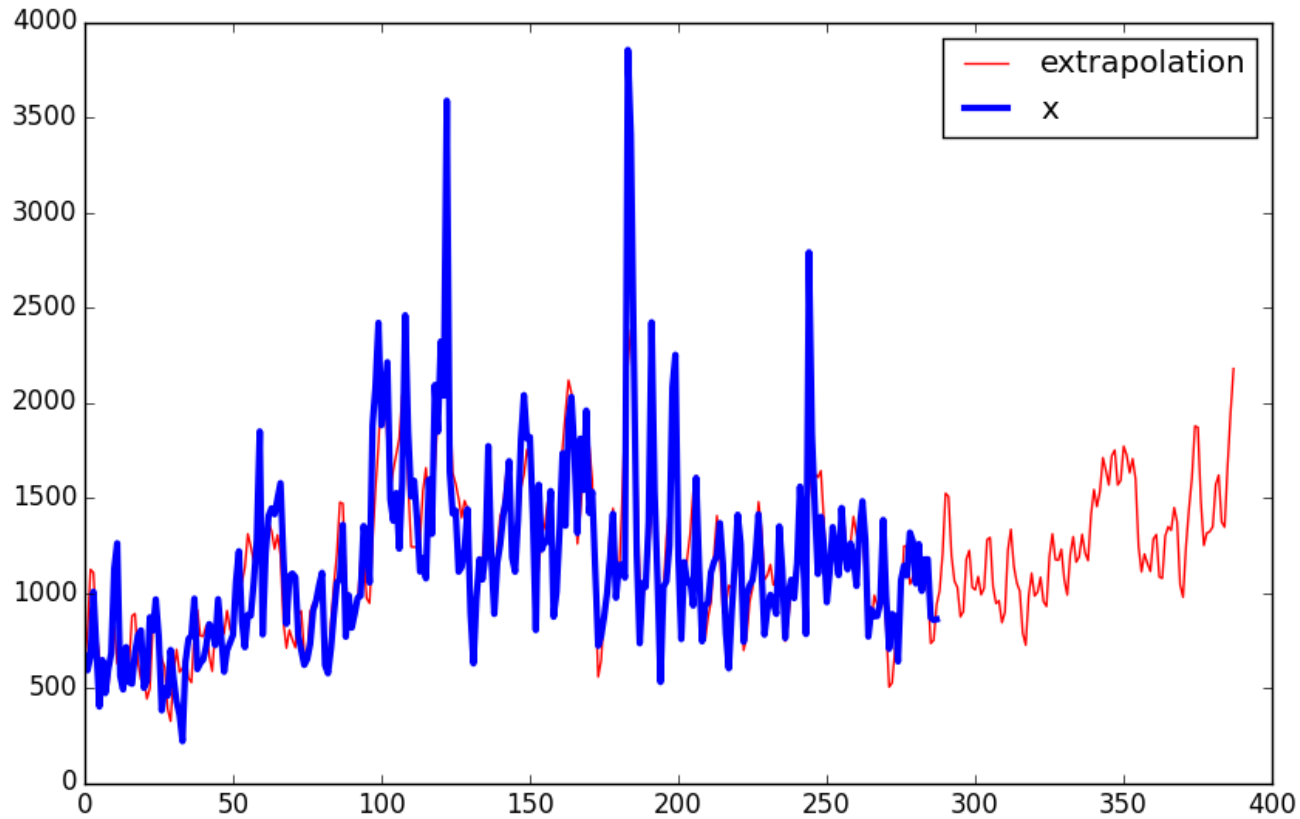
- Decompose a time series into a sum of many many sine or cosine functions
- The coefficients of these functions should have uncorrelated values
- Regression on sinusoids

Examples



1. What are the advantages?
2. Potential disadvantages?
3. When would this provide useful information?
4. When would this *not* provide useful information?

Fit can be surprisingly good



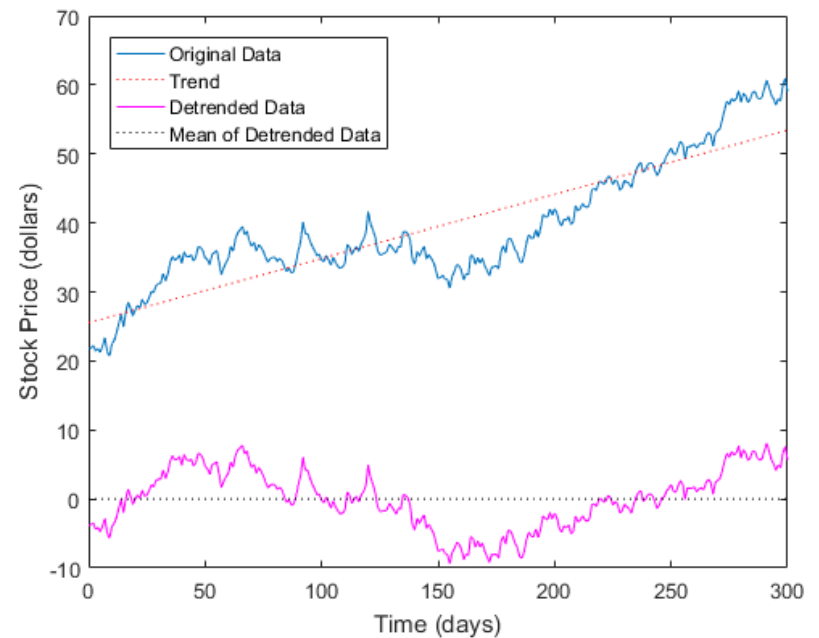
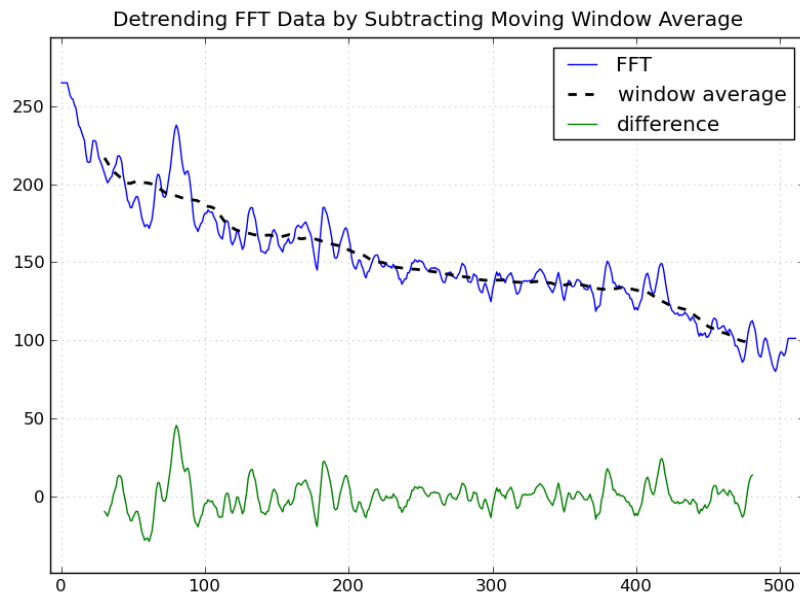
Pre-Prediction Munging

You need to remove the trend and seasonal elements before forecasting

- Most (interesting) data in the real world will show
 - Trends
 - Seasonality
- Most models require data that shows neither of these properties to say something interesting

De-trend your data

Use local smoothing or a linear regression

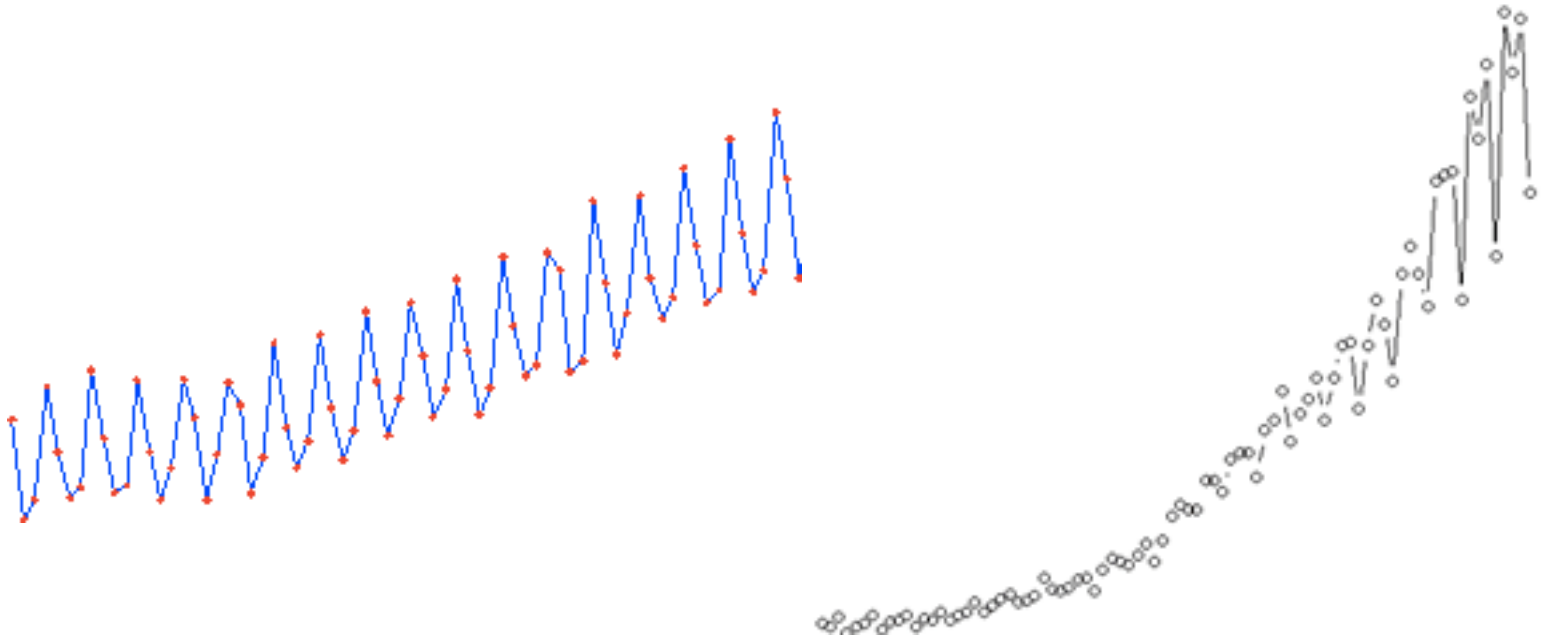


Remove seasonality

Look at power spectrum **after** removing trend data

- Simplest: average de-trended values for specific season
- More common: use 'loess' method ('locally weighted scatterplot smoothing')
 - Window of specified width is placed over the data
 - A weighted regression line or curve is fitted to the data, with points closest to center of curve having greatest weight
 - Weighting is reduced on points farthest from regression line/curve and calculation is rerun several times.
 - This yields 1 point on loess curve
 - Helps reduce impact of outlier points

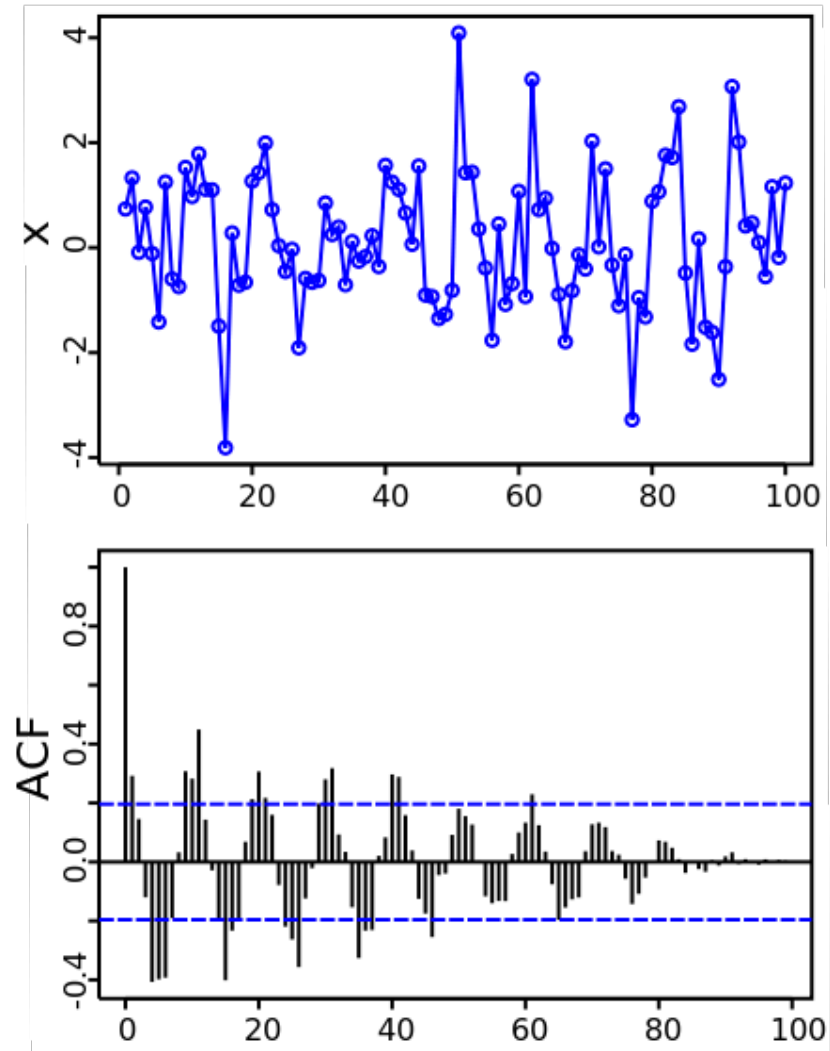
Additive vs. multiplicative seasonality



Prediction

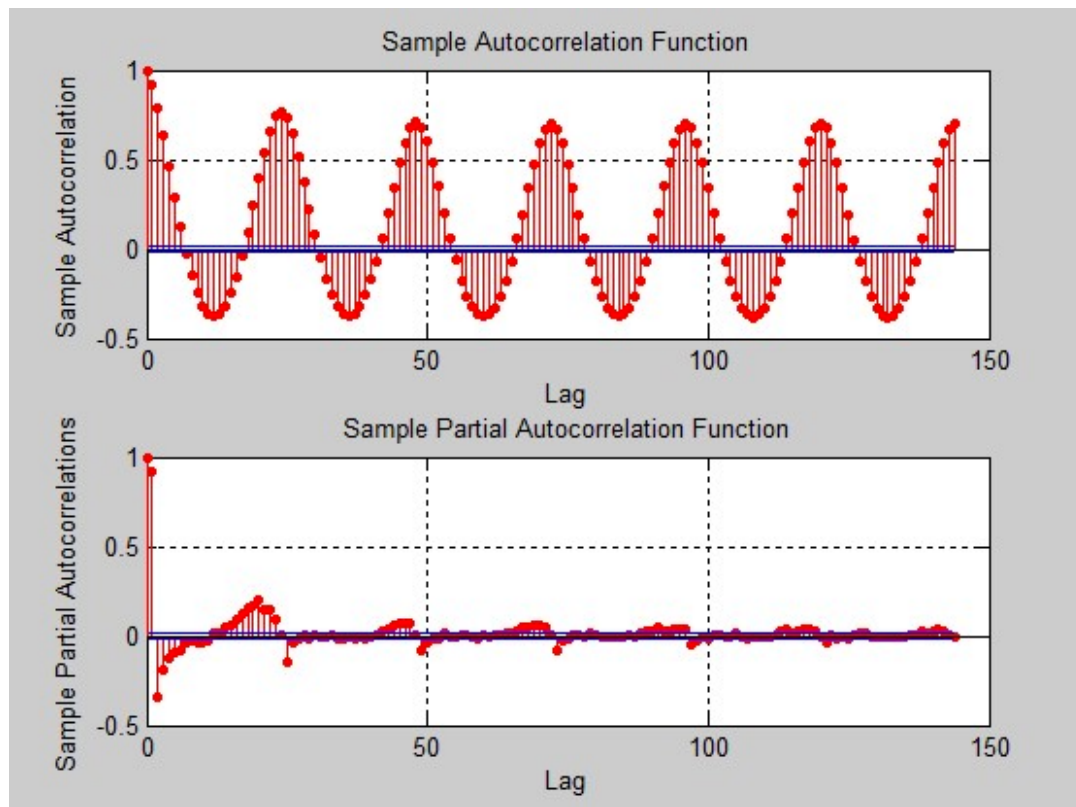
Autocorrelation function

- Used to help identify possible structures of time series data
- Gives a sense of how different points in time relate to each other in a way explained by temporal distance



Partial autocorrelation function

- “gives the partial correlation of a time series with its own lagged values, controlling for the values of time series at all shorter lags”
- Why would this be useful?



ARIMA model (a.k.a. Box-Jenkins)

- AR = autoregressive terms
- I = differencing
- MA = moving average
- Hence specified as (autoregressive terms, differencing terms, moving average terms)

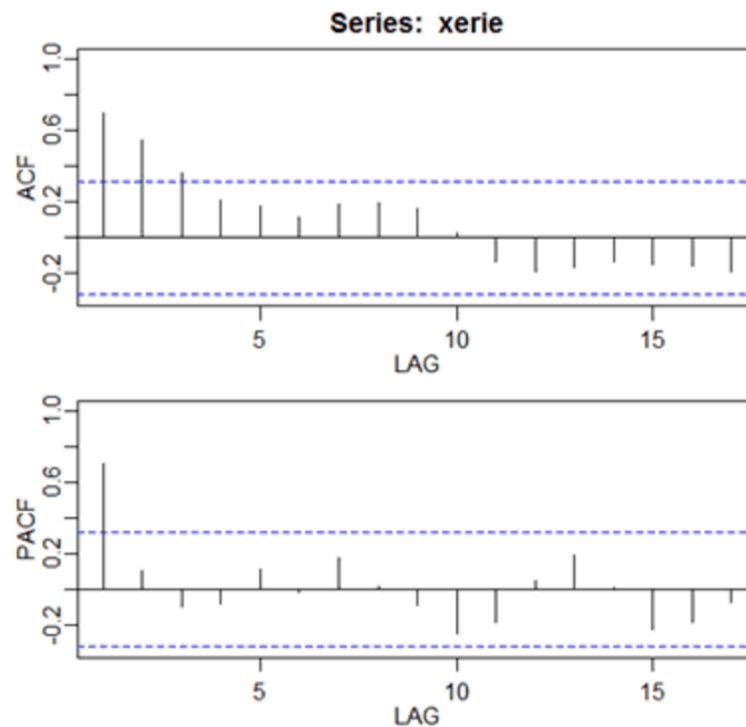
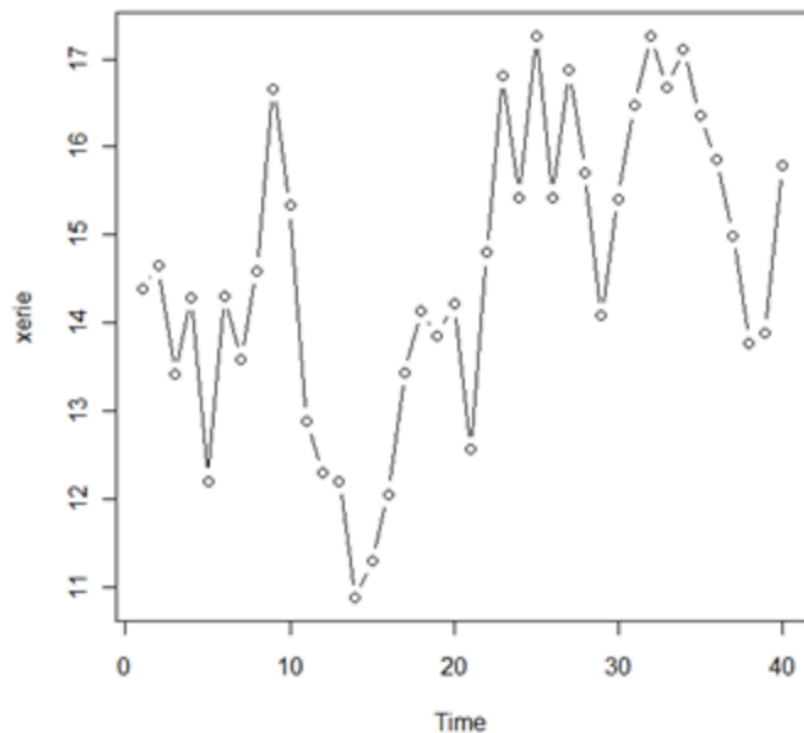
ARIMA mode: 'the most general class of models for forecasting a time series which can be **made to be 'stationary'**

- Statistical properties (mean, variance) constant over time
- 'its short-term random time patterns always look the same in a statistical sense'
- Autocorrelation function & power spectrum remain constant over time
- Ok to do non-linear transformations to get there
- ARIMA model can be viewed as a combination of signal and noise
- Extrapolate the signal to obtain forecasts

Applying the appropriate ARIMA model

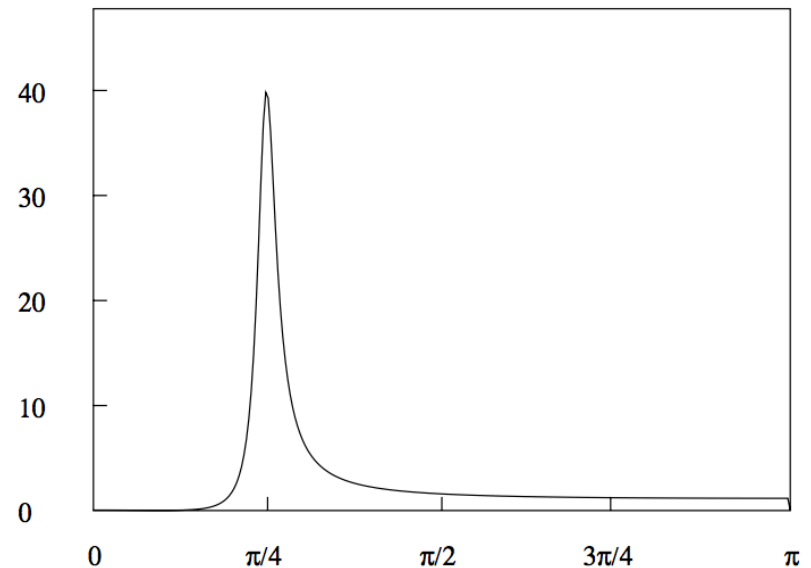
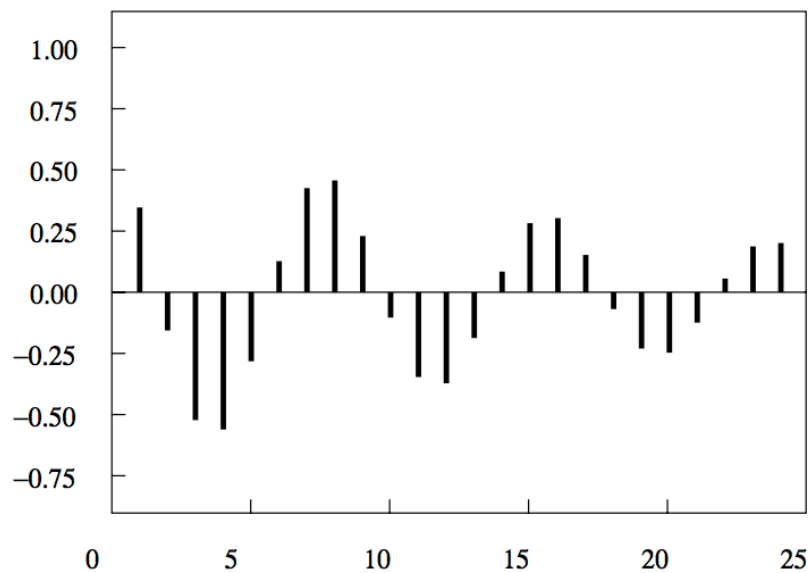
- Need to determine what ARIMA model to use
- Use plot of the data, the ACF, and the PACF
- With the plot of the data: look for trend (linear or otherwise) & determine whether to transform data
- Most software will use a maximum likelihood estimation to determine appropriate ARIMA parameters

Example – Annual Lake Erie depth



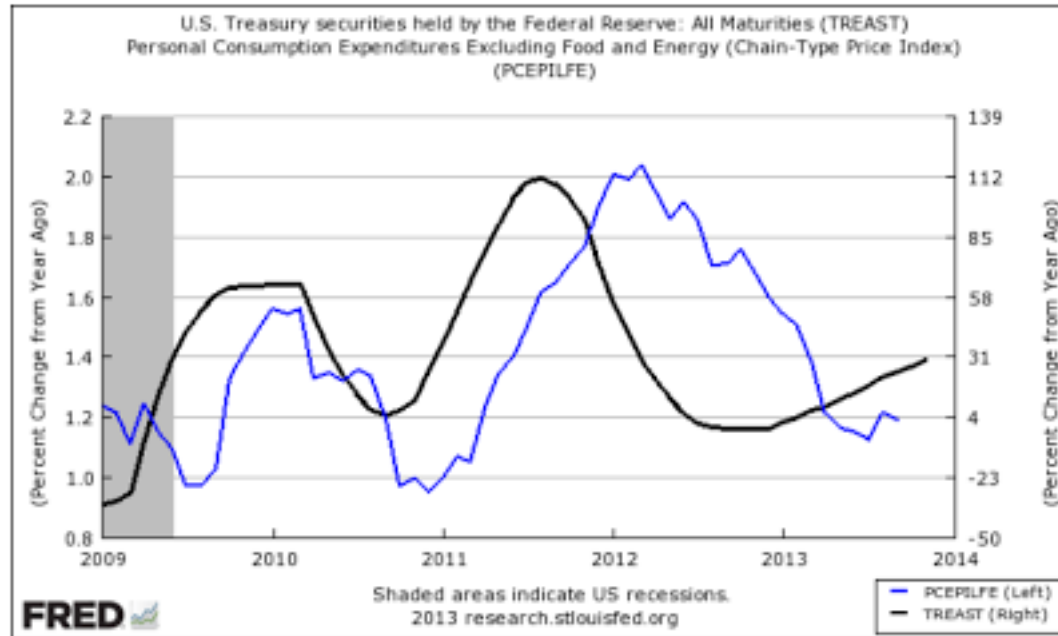
Spectral analysis is helpful for ARIMA modeling

Here we see an ARMA(2,2) process and its spectral decomposition



Learn more...

Vector autoregression works similarly for cases of multivariate time series



$$x_{t,1} = \alpha_1 + \phi_{11}x_{t-1,1} + \phi_{12}x_{t-1,2} + \phi_{13}x_{t-1,3} + w_{t,1}$$

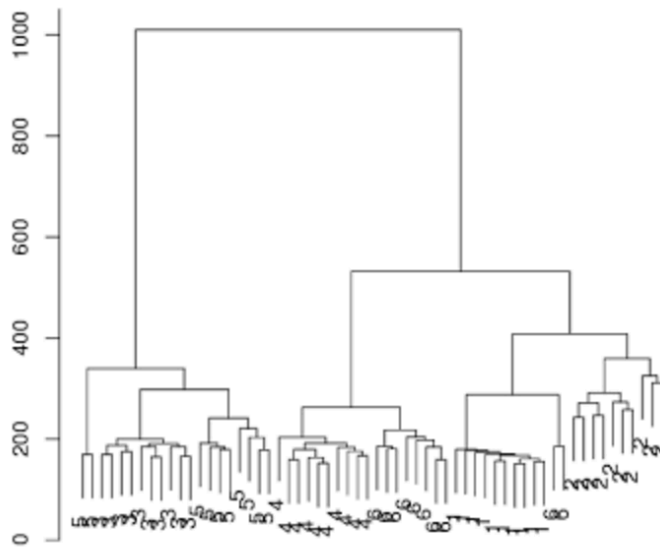
$$x_{t,2} = \alpha_2 + \phi_{21}x_{t-1,1} + \phi_{22}x_{t-1,2} + \phi_{23}x_{t-1,3} + w_{t,2}$$

$$x_{t,3} = \alpha_3 + \phi_{31}x_{t-1,1} + \phi_{32}x_{t-1,2} + \phi_{33}x_{t-1,3} + w_{t,3}$$

Clustering & Classification

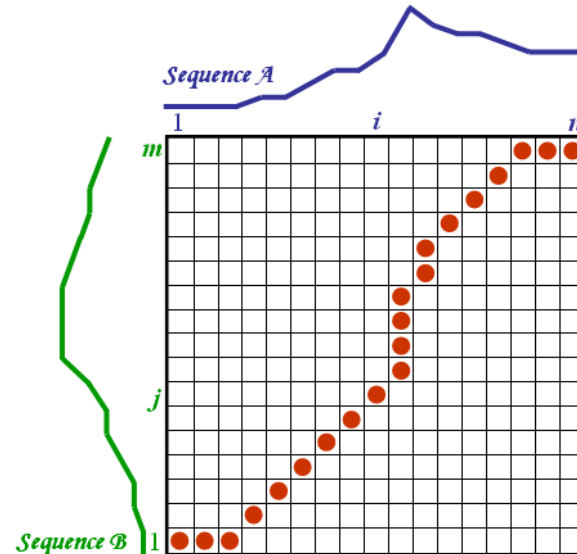
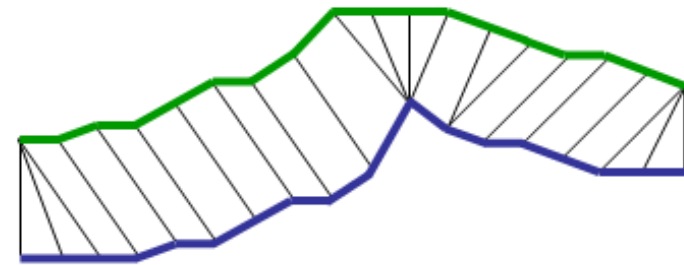
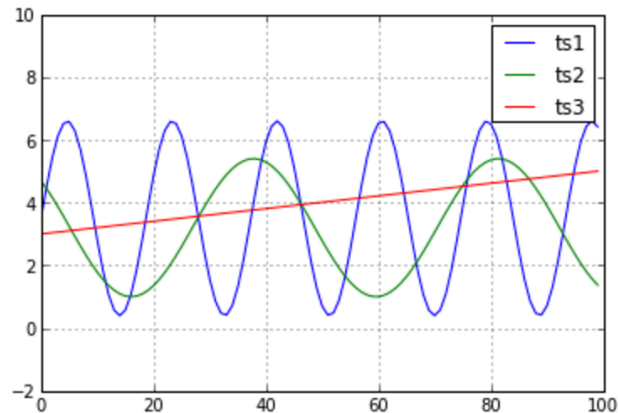
(yet another route to prediction)

Clustering



Choose an appropriate clustering technique based on what you know about the distance measure you use

Need to think carefully about distance metric



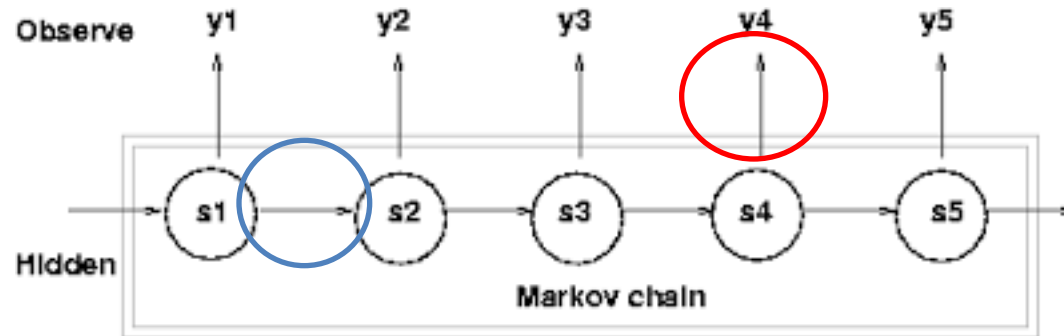
DTW-based Classification

DTW-based decision trees

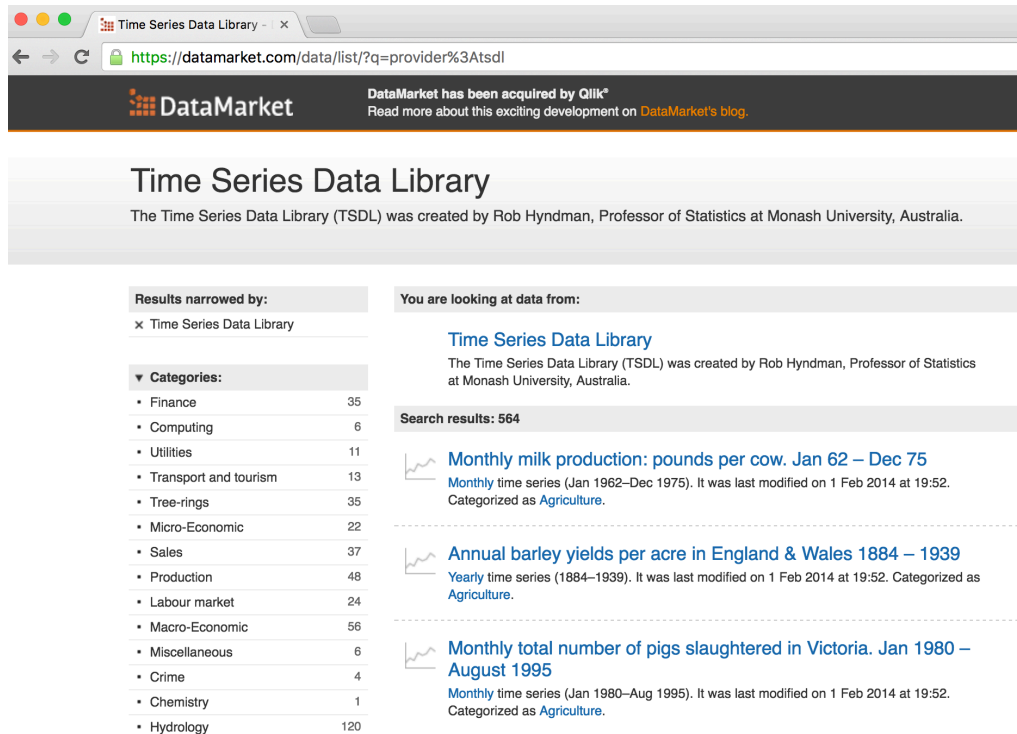
DTW-based nearest neighbor

Learn more...

Hidden Markov Models are another way of thinking about time series classification.



Get More Practice



The screenshot shows a web browser window with the URL <https://datamarket.com/data/list/?q=provider%3Atsdl>. The page title is "Time Series Data Library". Below the title, it states: "The Time Series Data Library (TSDL) was created by Rob Hyndman, Professor of Statistics at Monash University, Australia." The page is divided into two main sections. On the left, under "Results narrowed by:", there is a list of categories with their respective counts: Finance (35), Computing (6), Utilities (11), Transport and tourism (13), Tree-rings (35), Micro-Economic (22), Sales (37), Production (48), Labour market (24), Macro-Economic (56), Miscellaneous (6), Crime (4), Chemistry (1), and Hydrology (120). On the right, under "You are looking at data from:", there is a link to "Time Series Data Library" and a description: "The Time Series Data Library (TSDL) was created by Rob Hyndman, Professor of Statistics at Monash University, Australia." Below this, there is a section for "Search results: 564" which lists three time series: "Monthly milk production: pounds per cow. Jan 62 – Dec 75", "Annual barley yields per acre in England & Wales 1884 – 1939", and "Monthly total number of pigs slaughtered in Victoria. Jan 1980 – August 1995". Each entry includes a small line graph icon and a brief description of the data.

Time Series Data Library

The Time Series Data Library (TSDL) was created by Rob Hyndman, Professor of Statistics at Monash University, Australia.

Results narrowed by:

- Time Series Data Library

Categories:

- Finance 35
- Computing 6
- Utilities 11
- Transport and tourism 13
- Tree-rings 35
- Micro-Economic 22
- Sales 37
- Production 48
- Labour market 24
- Macro-Economic 56
- Miscellaneous 6
- Crime 4
- Chemistry 1
- Hydrology 120

You are looking at data from:

[Time Series Data Library](#)

The Time Series Data Library (TSDL) was created by Rob Hyndman, Professor of Statistics at Monash University, Australia.

Search results: 564

[Monthly milk production: pounds per cow. Jan 62 – Dec 75](#)

Monthly time series (Jan 1962–Dec 1975). It was last modified on 1 Feb 2014 at 19:52. Categorized as [Agriculture](#).

[Annual barley yields per acre in England & Wales 1884 – 1939](#)

Yearly time series (1884–1939). It was last modified on 1 Feb 2014 at 19:52. Categorized as [Agriculture](#).

[Monthly total number of pigs slaughtered in Victoria. Jan 1980 – August 1995](#)

Monthly time series (Jan 1980–Aug 1995). It was last modified on 1 Feb 2014 at 19:52. Categorized as [Agriculture](#).

www.cs.ucr.edu/~eamonn/time_series_data/



UCR Time Series Classification Archive