

TIME SERIES

2020.05.11
Mags Blust



Unsupervised Learning

DIMENSIONALITY REDUCTION

PCA
t-SNE
Factor Analysis

CLUSTERING

AD, NLP may be solved with clustering.

Algorithms

KMeans
DBScan
Hierarchical

MACHINE LEARNING

Reinforcement Learning

AD, NLP, RS may be solved with reinforcement learning

Collaborative Filtering
Content Filtering
Neural Networks (DL)

Supervised Learning

CLASSIFICATION

- Discrete, Categorical target

TSA, NLP, AD may be classification problems

REGRESSION

- Continuous Target

TSA, AD maybe regression probs.

Linear Regression
GLM
Polynomial Reg.
SVR
Decision Tree Regr.
Neural Networks

Logistic Regression
Decision Tree
Random Forest
SVC
KNN
Neural Networks

SUB-METHODOLOGIES

TSA

time series analysis

AD: anomaly detection

NLP: natural language Processing

RS: recommendation System

TIME SERIES analysis

Time Series Forecasting in Minutes: https://www.youtube.com/watch?v=wGUV_XqchbE

"TIME is not a line but a series of Now Points" ~Taisen Deshimaru

- finding patterns in temporal data, and making predictions
- a sub-methodology to Regression / classification- (predominantly) regression

Why do we treat it differently than "normal" regression?

general Linear Regression

X=features									target
x ₁	x ₂	x ₃	x ₄	x ₅	x ₆	x ₇	x ₈	x ₉	y

Time Series

X=features=time!								target
t ₋₇	t ₋₆	t ₋₅	t ₋₄	t ₋₃	t ₋₂	t ₋₁	t ₀	

ASSUMPTION: features are independent of each other !

BUT : these features (consecutive date/time stamps) are, by their nature, dependent on each other !

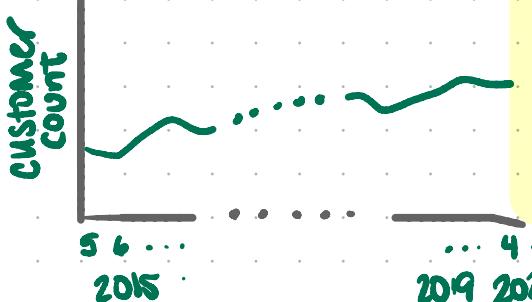
TIME SERIES

forecast/predict # of new customers next month using historical monthly new customer count

- Features = each historical Month
- target = next Month

1 dimension

Date	Customer Count
2015-05	
2015-04	
:	:
2020-04	
2020-05	



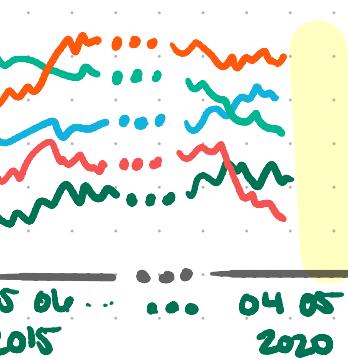
Customer count

5 6 ... 2015 ... 4 5 2019 2020

5 dimensions

Product	5/15	6/15	...	4/20	5/20
savings acct					
checking acct					
home insurance					
auto insurance					
renters ins.					

Customer Count



1 Model

NOT TIME SERIES

Predict # of new customers each month using Sales/Marketing activity.

- observation = Month/Gear
- features (e.g.) = prev. Month campaign count, 30 d. Δ in sales staff, # New Prod. I features

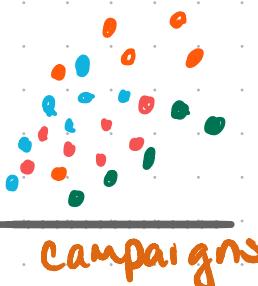
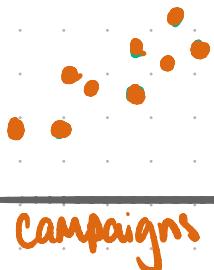
- target = customer count

FEATURES → Y

Month-Year	# of campaigns	Δ in sales staff	...	Customer count
5-15				
6-15				
:				
3-20				
4-20				

Model patterns from previous months to be able to give a 30 d. forecast of customer count.

Customer Count



To predict by product, build multiple models, 1 per product

Savings acct
Checking acct
Home insurance
Auto insurance
Renters insurance

Time Series Vocabulary

RESAMPLING (dates) Changing frequency of data points.

TREND: Long term progression (increasing or decreasing)



SEASONALITY: Series is influenced by seasonal factors

e.g. Month of year, day of week.

Always a fixed, known period. \Rightarrow

Seasonal series == PERIODIC series



date	New customers
2020-04-01	50
2020-04-02	125
:	:
2020-04-30	75
2020-05-01	53

date	New customers
2020-04	50+125+...+75
2020-05	53+...

Resample Monthly & SUM

CYCLIC: fluctuations that are NOT of a fixed period. Duration of fluctuations > 2 yrs
e.g. housing market.

HETEROSKEDASTICITY: Changes in Variance over TIME

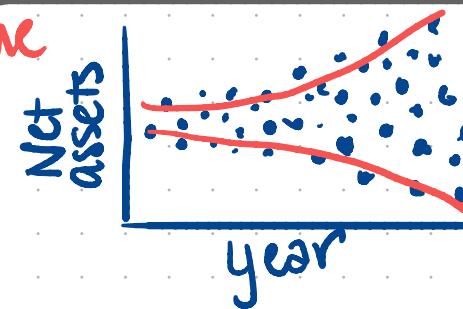
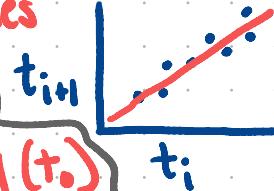
AUTOCORRELATION: "Regression of Self"

Used to detect non-randomness in data -

It is a correlation coefficient, but instead of between 2 different variables, it is between the values of the same variable at 2 different times

LAG VARIABLE: Previous time step

yesterday ($t-1$) is a lag var to today (t_0)



"tomorrow's" value is dependent on "today's" value.

New Skills

Acquisition: gather data using a Rest API

Prep: Working with dates - `resample()`
`asfreq()`

filling missing values - `fillna()` `ffill`, `pad`, `bfill`

Explore: Splitting TS data into train / test - `sklearn.model_selection.TimeSeriesSplit`
- using date cutoff

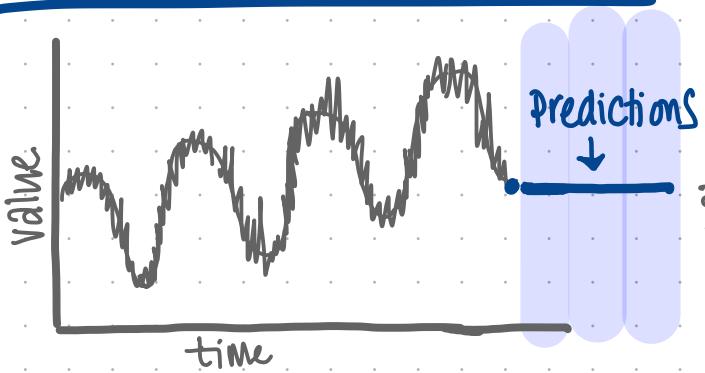
Viz time series data:
Plot simple aggregate (`resample().mean().plot()`)
Rolling aggregate (`rolling().mean().plot()`)
Pd. over Pd. diff. (`resample().mean().diff().plot()`)
Customize datetime axis `matplotlib.dates`
Explore 'Seasonality'
Merge multiple frequencies into same plot
Time Series Decomposition
Analyzing Lag

Model: Forecast / Predict Methods include

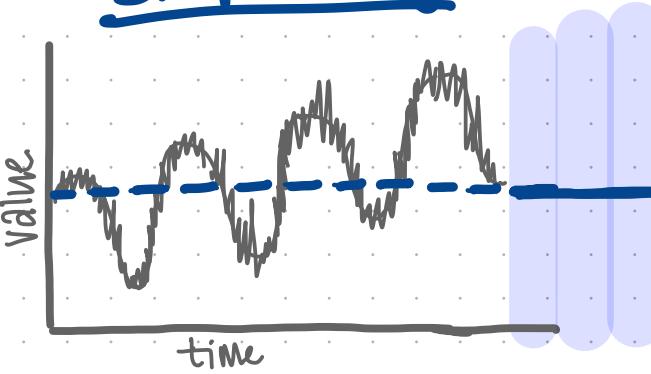
- last observed value
- simple avg.
- moving/rolling avg
- Holt's linear trend
- previous cycle
- FB Prophet Module.

Forecasting, Predicting, Modeling Time Series Data

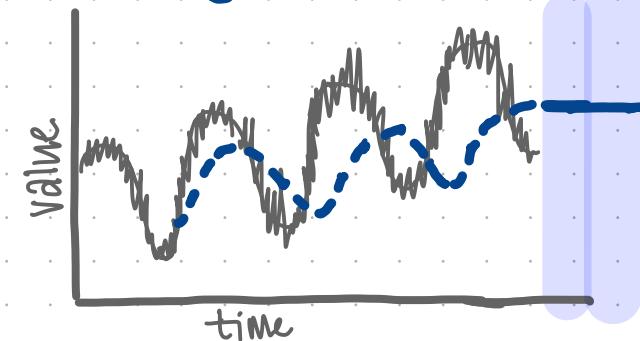
Last Observed Value



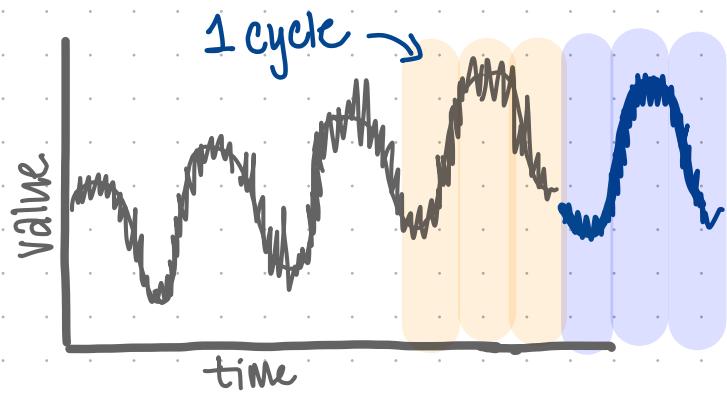
Simple Avg



Moving / Rolling Avg.

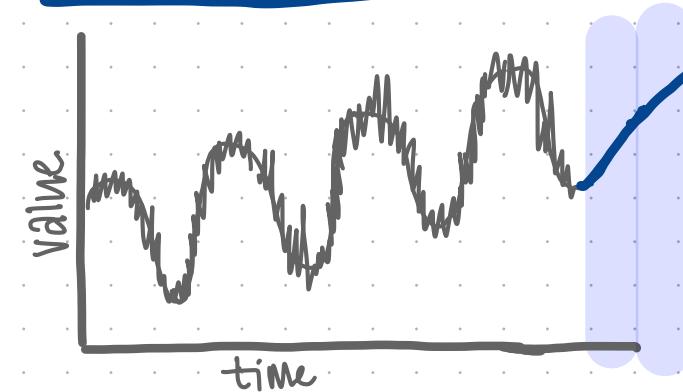


Previous Cycle



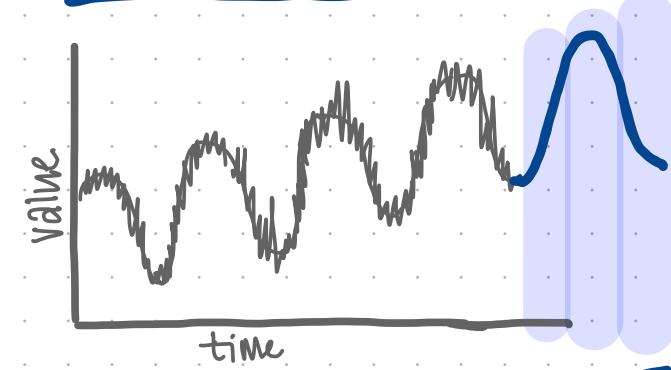
[define a cycle, predict the next cycle to be the values of the previous cycle]

Holt's Linear Trend



[Exponential Smoothing applied to both the average & the trend (slope)]

FB Prophet Model



[Non-linear trends fit with yearly, weekly, daily seasonality + holiday effects]