

Machine Learning

SUPERVISED

CLASSIFICATION REGRESSION

TARGET:

discrete
categorical

continuous
numeric ordered

SUB-METHODOLOGIES:

time series
Anomaly detection
NLP

time series
anomaly detection

ALGORITHMS:

Logistic Regression
SVM - Classifier
Decision Tree
Random Forest
KNN
Neural Network
Naïve Bayes

OLS (linear regression)
GLM (Hawdeie regression)
Polynomial Regression
Support Vector Regressor
Decision Tree Regressor
Neural Networks
LASSO, LARS
Holt-Winters
Auto Regression

TIME Series Module

UNSUPERVISED

CLUSTERING DIMENSIONALITY REDUCTION

SUB-METHODOLOGIES:

Anomaly Detection
NLP

ALGORITHMS:

K-Means
DBSCAN
Hierarchical

PCA
t-SNE
factor Analysis

REINFORCEMENT

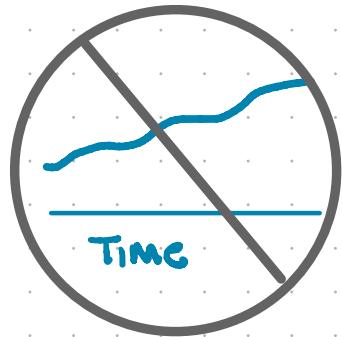
SUB-METHODOLOGIES

Anomaly Detection
NLP
Recommender Systems

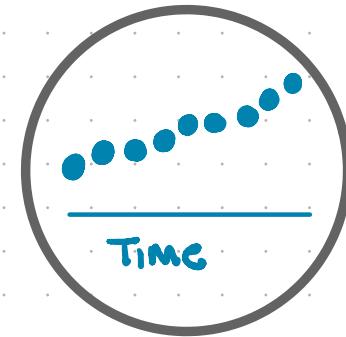
ALGORITHMS:

Collaborative Filtering
Content Filtering
Neural Networks

TIME SERIES ANALYSIS



"Time is not a line but a series
of 'now' points"
Taisen Deshimaru



TSA is finding patterns in temporal data
& making predictions.

TSA is a sub-methodology to Regression &
Classification. ∴ It is a form of supervised
Machine learning

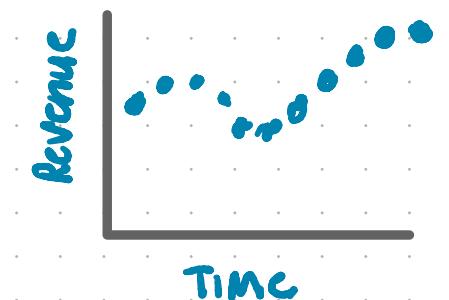
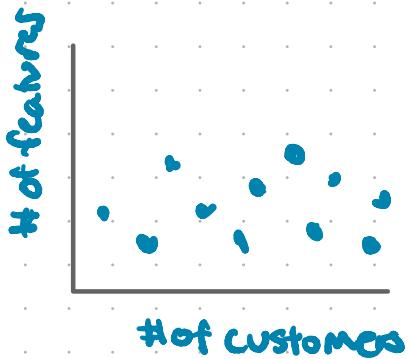
TSA - a special case of regression - WHY?

because... **ASSUMPTIONS**

TIME Series data does **NOT** meet required assumptions for linear regression.

Specifically... **COLLINEARITY**

The features in TSA are date/times.
By their nature, they are dependent on each other.



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
1	2	3	4	5	6	7	8	9	10

ASSUMPTION: features are independent of each other!

Time Series							
X=features = time!							
t-7	t-6	t-5	t-4	t-3	t-2	t-1	t ₀
last week							today
							rest.

X = dependent

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

TO TSA or not to TSA

...THAT IS THE QUESTION

SCENARIO: Predict number of new customers we will have (customer growth)
OPTIONS

TIME SERIES REGRESSION

use previous outcomes as predictors
to predict future outcomes

	Customer growth		
	TRAIN	VALIDATE	TEST
Jan 20	100	100	100
Feb 20	90	90	90
Mar 20	95	95	95
Apr 20	80	80	80
May 20	105	105	105
Jun 20	110	110	110
Jul 20	110	110	110
Aug 20		90	90
Sep 20			70
Oct 20			

TRAIN

LABELED OUTCOME

predictors

VALIDATE

LABELED OUTCOME

predictors

TEST

LABELED OUTCOME

predictors

METHODOLOGY

ORDINARY REGRESSION

use Attributes to predict outcome

	customer growth	Sales staff	Mktg \$	# campaigns	# new products
Jan 20	100				
Feb 20	90				
Mar 20	95				
Apr 20	80				
May 20	105				
Jun 20	110				
Jul 20	110				
Aug 20	90				
Sept 20	70				
Oct 20					

* Attributes = predictors = features = independent vars.
* outcome = target = dependent variable

Time Series Vocabulary

RESAMPLING (dates) Changing frequency of data points.

TREND: Long term progression (increasing or decreasing)



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

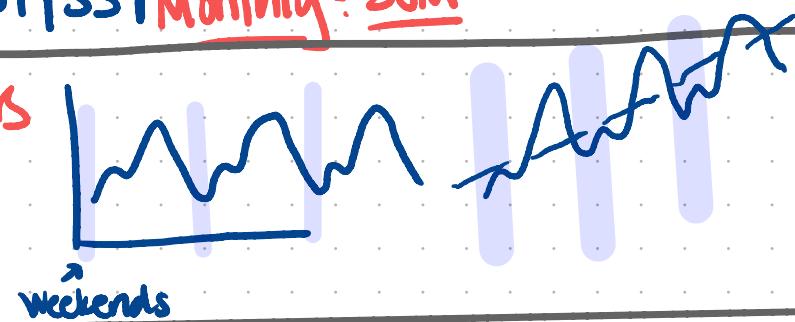
Resample Monthly & SUM

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



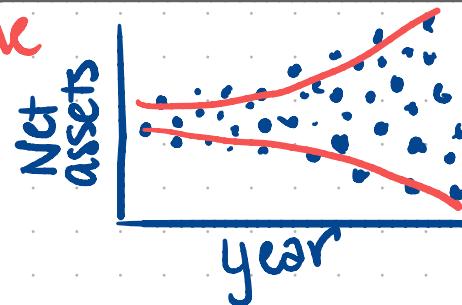
CYCLIC: fluctuations that are NOT of a fixed period. Duration of fluctuations > 2 yrs
e.g. housing market.

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



COLLINEARITY
when independent variables are highly correlated

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.

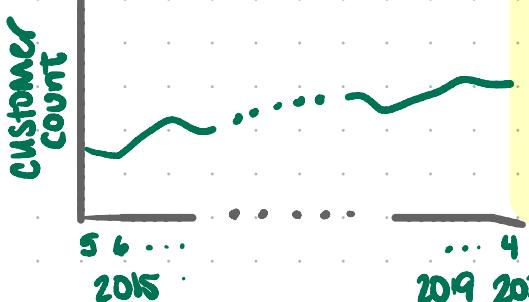
TIME SERIES

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	

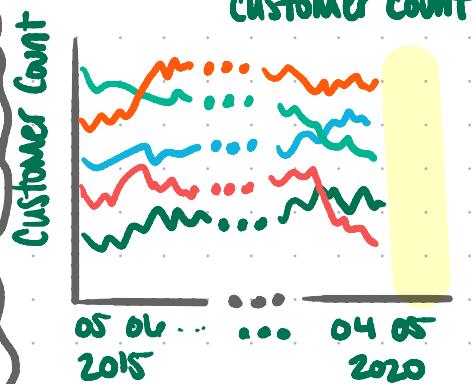


Customer count

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.					



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

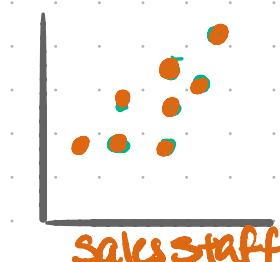
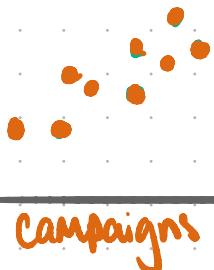
- target = customer count

FEATURES → Y

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

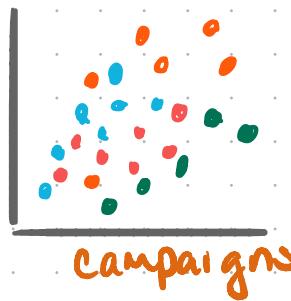
Model patterns from previous months to be able to give a 30 d. forecast of customer count.

Customer count



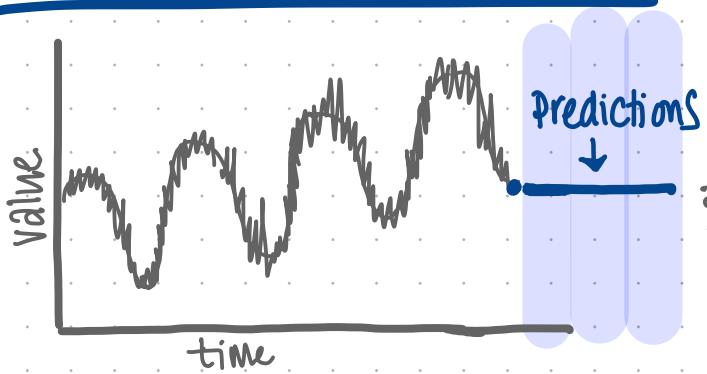
To predict by product, build multiple models, 1 per product

Savings acct
Checking acct
Home insurance
Auto insurance
Renters insurance

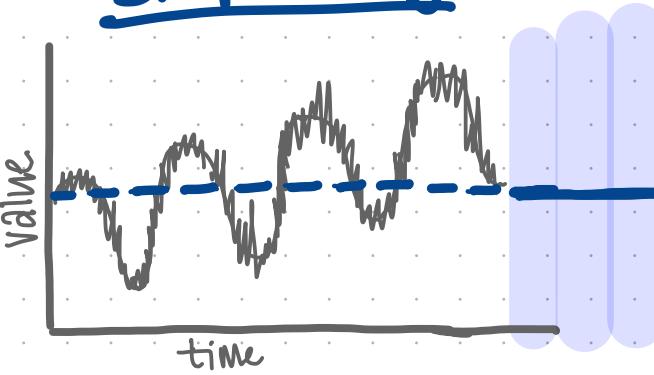


Forecasting, Predicting, Modeling Time Series Data

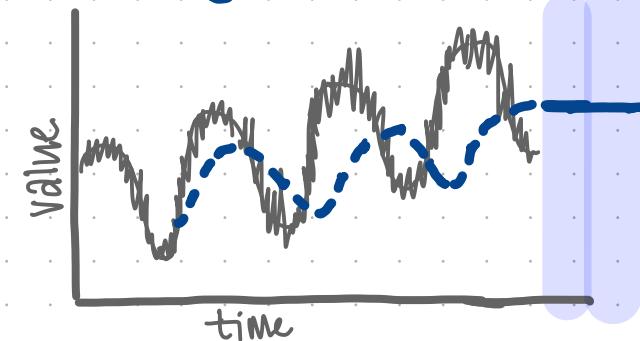
Last Observed Value



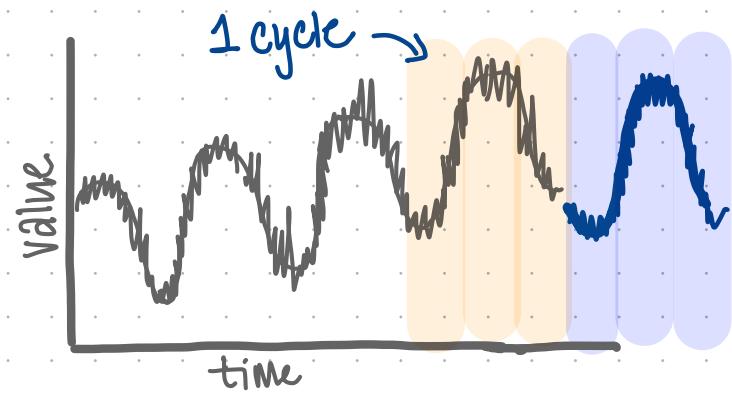
Simple Avg



Moving / Rolling Avg.

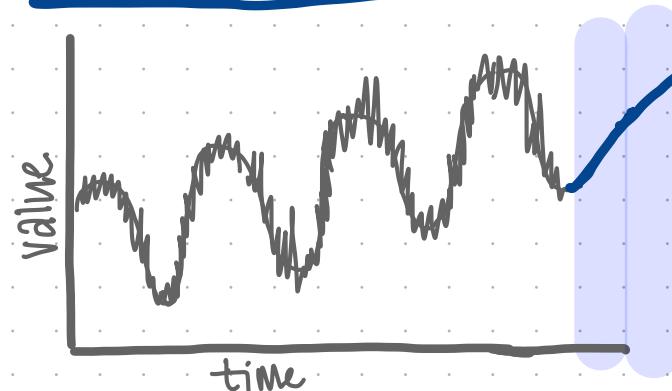


Previous Cycle



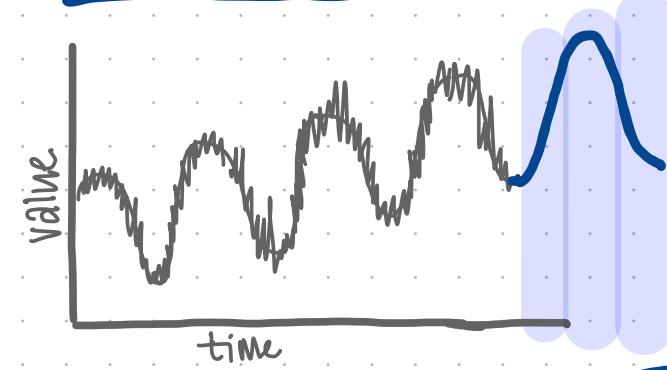
[define a cycle, predict the next cycle to be the values of the previous cycle]

Holt's Linear Trend



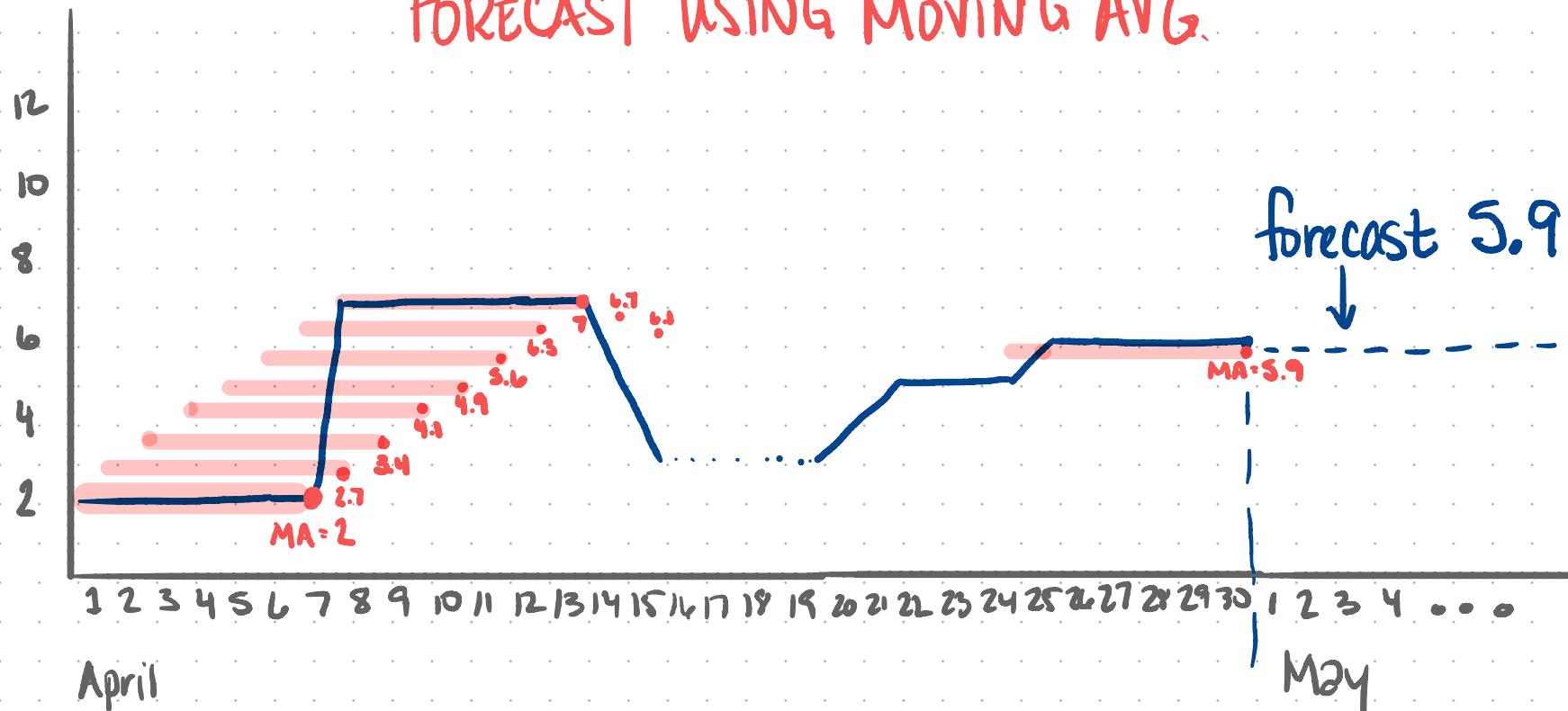
[Exponential Smoothing applied to both the average & the trend (slope)]

FB Prophet Model



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

FORECAST USING MOVING AVG.



7-day Moving/Rolling Average

Time Series Cross Validation - Prophet

