

Utility: Predictive Analytics for Business

Problem Solving with Analytics

Common Industry Standard Process for Data Mining (CRISP-DM)

1. business understanding - what is decision, what info needed, type of analysis
2. data understanding - what data needed, available, type
3. data preparation - gather, cleanse, format, clean, sample
4. analysis / modeling - methodology, identify important factors
5. validation
6. presentation / visualization
 - ↳ custom to audience
 - walk audience through audience
 - tell story
 - reference data sources

Methodology Map:

Business Problem					
Predict Outcome			Data Analysis		
Data Rich		Data Poor		Geo Spatial	
Numeric		Classification		Segmentation	
Continuum	Time Based	Binary	Non-Binary	Aggregation	
- Linear Regression	- ARIMA	- Logistic Regression	- Forest Model	Descriptive	
- Decision Tree	- ETS	- Decision Tree	- Boosted Model		
- Forest Model					
- Boosted Model					

Linear Regression: $y = mx + b$

$$M = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$\text{coefficient of determination, } R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

$$\text{sum of squared total, } SST = \sum (y_i - \bar{y})^2$$

$$\text{sum of squared regression, } SSR = \sum (\hat{y}_i - \bar{y})^2$$

$$\text{sum of squared error, } SSE = \sum (y_i - \hat{y}_i)^2$$

multiple linear regression,

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots$$

$$= b_0 + \sum b_i x_i$$

Data Wrangling: Understanding Data

3 types of data structure

- structured - organized into rows / columns

- unstructured - no structure

- semi-structured - some structure, e.g. log, require parsing.

Data sources

- computer files - e.g. XLS, XML, CSV

- database - e.g. mongoDB, MySQL

- web-based - salesforce, twitter, ..

Data Types

- string - alphanumeric

- numeric - integer, float,

- date/time - date, time

- boolean - True/False

- special objects - images, mp3, pdf, audio files, video files, ...

Data Issues:

- Dirty data
- missing
 - duplicate
 - unusable format - not parsed properly
 - extra character / missing characters
 - incorrect data

Parsing Data - CSV

- delimiter - ',', ';' } text to column

Extra characters - " ", \$, ^M

- Formula - Replace

Left - length of character

Duplicate - Unique

Sum (groupby, count) + Filter
" " + Sort

Missing - some algos won't work with missing data

adds bias

- delete using filters is Not Null

- explore using Field Summary

- impute - replace with mean, median, mode

- Summary + Append + Formula
(if, else)

or use Imputation

- others - multiple imputation - create model to predict

- full information maximum likelihood - factor missing values
into model

- number of missing vs significant of values

Outlier - incorrect data or abnormal but correct data

- box-whisker plot - $> 1.5 \times \text{interquartile range}$ beyond 1st, 3rd quartile
- build models w/o outlier + compare
- drop only if obvious outlier
- truncation - limit to max.

Data Formatting:

Transpose - Transpose tool

Aggregate - Sum tool

Cross tabulation - Join 2 tables

Data Blending:

Union - appends data

- Union tool

Join - join tables

- Join tool - left, join, right

Fuzzy matching - Jaro - measures characters in common

Levenshtein - counts number of edits to convert one string to another

- match 2 to 3 columns for duplicates, repeat for all permutations

Spatial Matching

~~points~~ points

- line - strings of latitude, longitude
- polygon - vertices of "

Spatial Blending

- Create Points, Trade area tools, Spatial match tool

Predictor Variables:

Matrix of correlation

1. Pearson correlation coefficient

↳ normalized covariance by standard deviation

↳ $\frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$ when applied to population

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \quad \text{when applied to sample}$$

2. Spearman's rank-order correlation

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

3. Hoeffding's independence test

Heat plot to visualize correlations - positive, negative.

↳ Association tool

Preparing data - EDA

↳ Field Summary tool - interesting - Browse tool

↳ Null/missing values - filter my (low % missing)
remove field (high % missing)
impute records (if important)

↳ Frequency tool - interesting - Browse tool

Classification Models

Binary classification - Logistic Regression

$$\text{multiple linear regression, } y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

$$\text{logistic regression, } \ln \left(\frac{P}{1-P} \right) = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

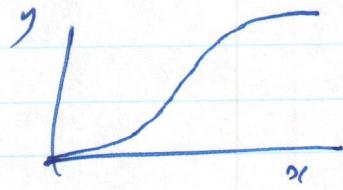
P = probability of outcome

$$P = \frac{e^{(b_0 + b_1 x_1 + \dots + b_n x_n)}}{1 + e^{(b_0 + b_1 x_1 + \dots + b_n x_n)}}$$

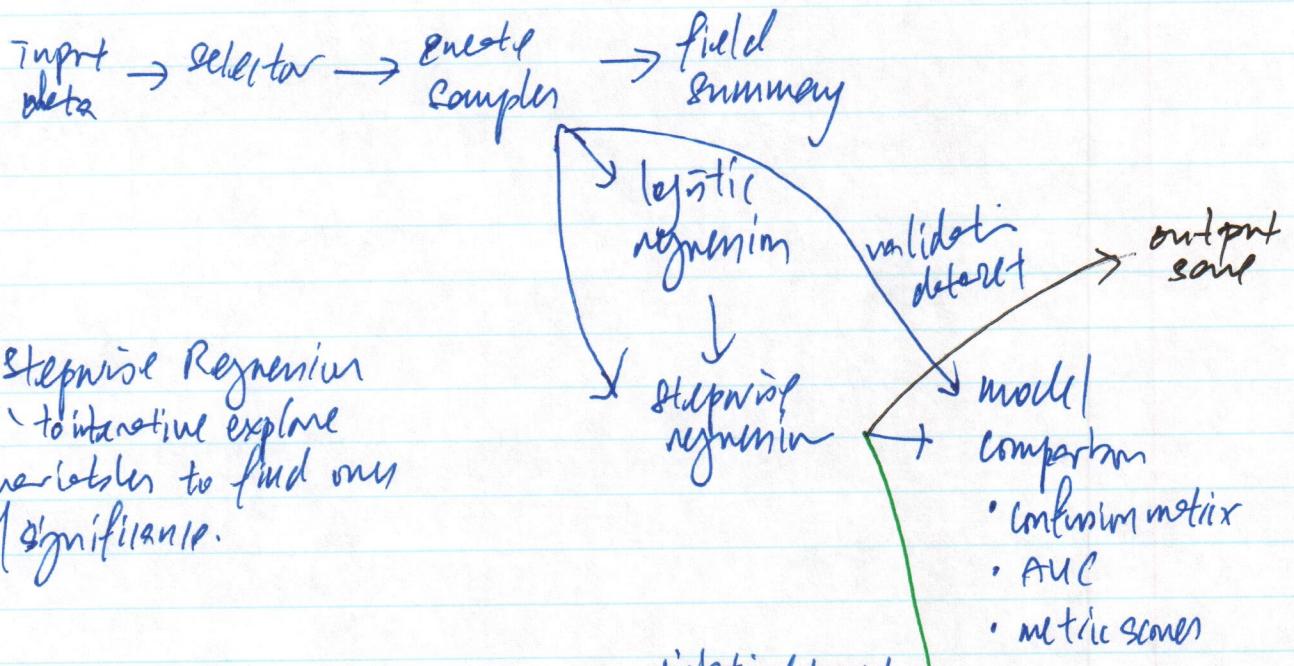
logit: $\eta(p) = \ln\left(\frac{P}{1-P}\right)$

probit: $\eta(p) = \Phi^{-1}(p)$

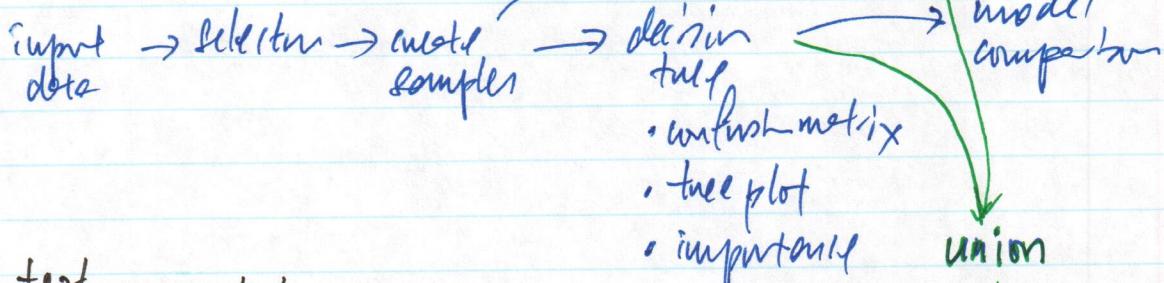
log-log: $\eta(p) = \ln(-\ln(1-p))$



Altayx workflow:



Decision Tree



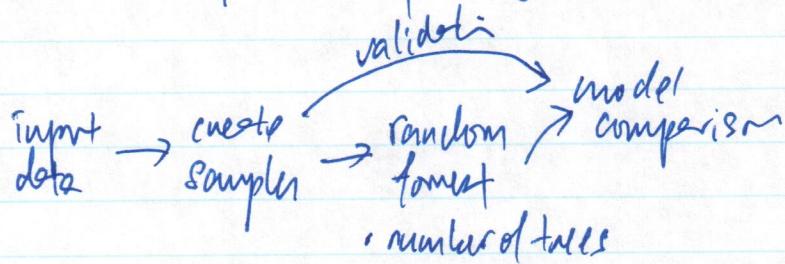
test data → selector

Input build model → score → browser

model comparison
• gain chart
• AUC
• precision-recall

Non-Binary Classification Model:

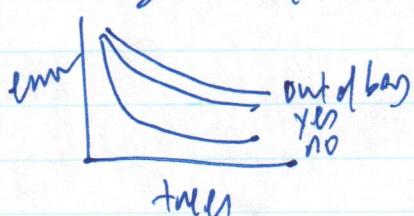
- Decision tree - tends to overfit training dataset
 - need to validate with validation test set
- Random forest - creates many trees to form ensemble of decision trees
 - each tree created by different randomly selected chunks of original data \rightarrow bootstrap
 - looks at overall as a whole to make prediction
 - overfitting gets averaged out by all trees.
 - majority vote over all terminal nodes of each tree.
probability averaged over all trees.



- Out-of-Bag Error Rate - indicates model performance with own-validation set in training dataset.

- Confusion matrix - true/false positive / negative

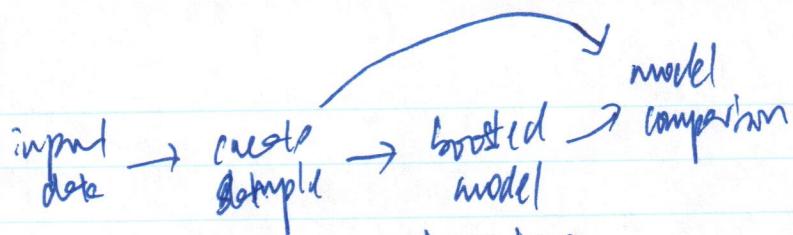
- Percentage Error for Trees



- number of trees needed to minimize error
i.e. convergence

- Feature importance by mean decrease in Gini

• Boosted Model - calculate errors from trees
modify tree to reduce error
repeats

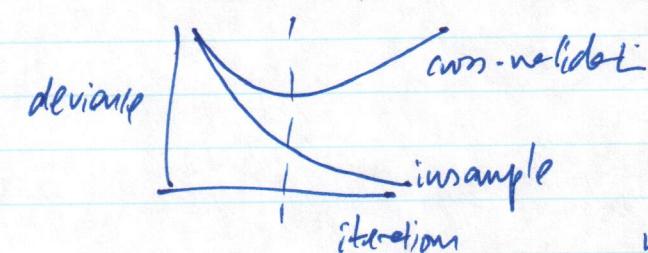


- target type
 - continuous,
 - binary

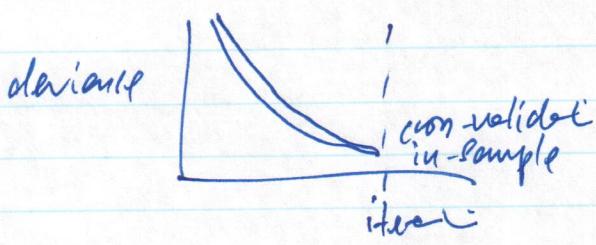
multinomial - 3 or more categories

- max trees
- cross-validation - k folds

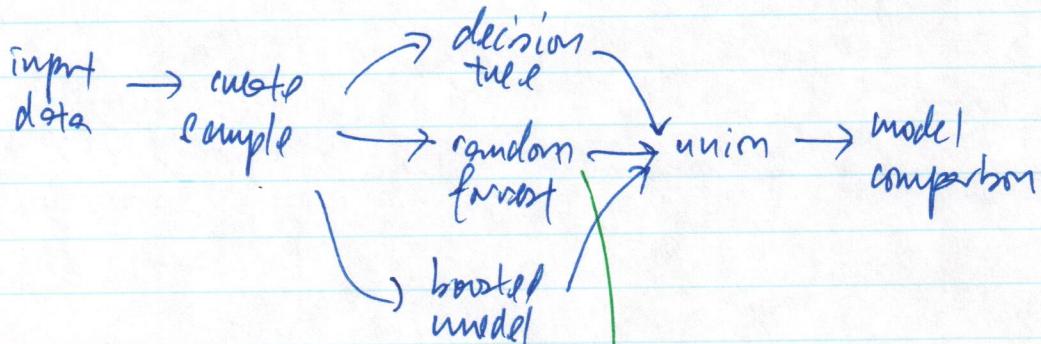
→ ~~import~~ feature importance plot - which is important
→ number of iterations assessment plot



- minimum
- a number of iteration needed
- large differences between in-sample & training + cross-validation indicate poor performance with validation data.



- similarity between cross-validate + in-sample, good predict with validation data. performance



test data

• binarize majority

• score → formula → summary

A/B Testing

Unit - person
store
product line
website

Treatment Group:

- collect units getting treatment
- treatment = change

Control Group:

- baseline comparison
- no treatment applied
- control units

Experimental Variable

- treatment change to see effect on target variable.
- independent

Control Variables

- ~~predictor~~ sales
- customer demographics
- location
- make sure treatment units + control units are similar to make comparison on experiment variables.

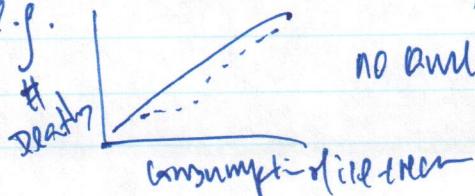
Steps to identify control variables

1. List of potential variables
2. Is data available
3. logical connection between control + target variable
4. Test ~~correlation~~ correlation between control + target variables
5. Test correlation between control variables.

Lurking variables / confounding variable

- highly correlated to target, controlling variable
- care to include controlling variable that nearly is not important.

- e.g.



no conf-effect, but lurking is temperature.

use logic to avoid.

Experiment Design

- Randomized Design

- ~ treatment, control units selected randomly on the fly.
- ~ used when hard to control variables, volume & velocity of data is high enough to avoid bias.
- ~ e.g. website, phone-based experiments.

- Matched Pairs

- ~ volume of observation is low
- ~ concern for bias is high
- ~ cost per observation is high
- ~ treatment unit matched to control unit
- ~ pair analysis, aggregate to single result.

Length of experiment

- minimum 1 full cycle of data to reduce chance of bias in responses.

Randomized Design Example:

- Treatment Control
- target variable - click on link - binary yes/no
- control variables - customer demographics
membership existing
- unit of diversion - how to assign units to control or treatment groups
- population - people without app
- duration } common control + treatment groups represent population
 - e.g. 1 week
unique visitor tracked by cookie or IP.
- size }
 - e.g. 2 servers - control
 - 1 server - treatment

test of means analysis

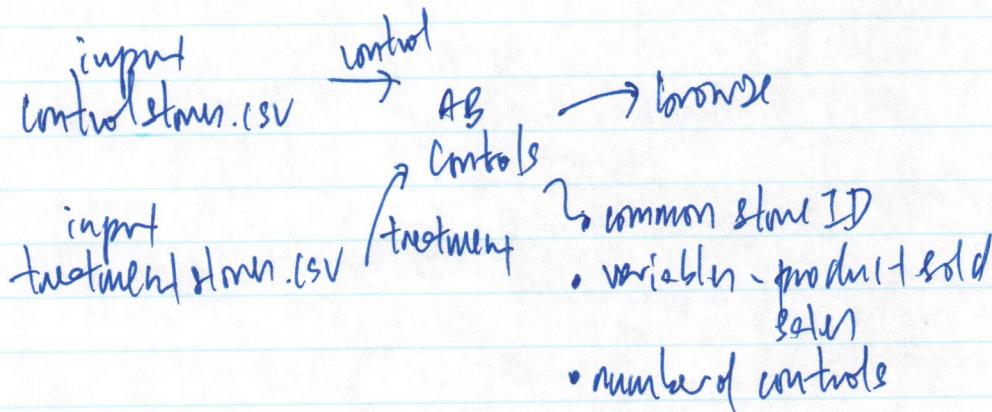
- test of whether mean of treatment, control groups are same
- $p = \text{likelihood difference between means} \Rightarrow \phi$
- $p < 0.05$ statistically significant.
- Test of Means tool

Matched Pair Design Example:

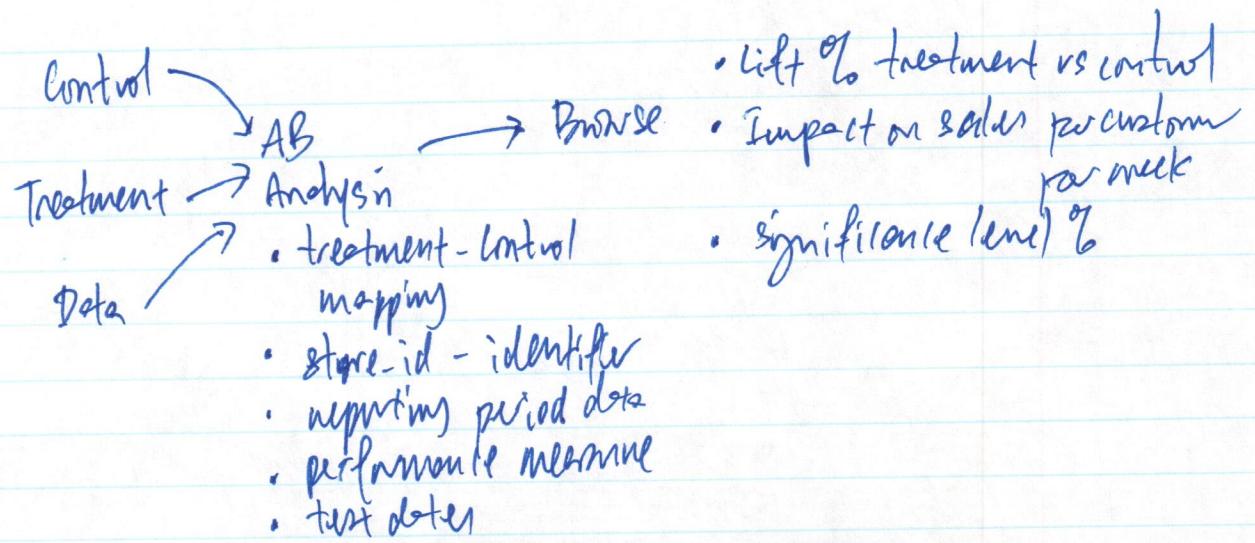
1. identify units for treatment units.
2. need > 10 treatments
3. randomly selected treatment units.

Control units

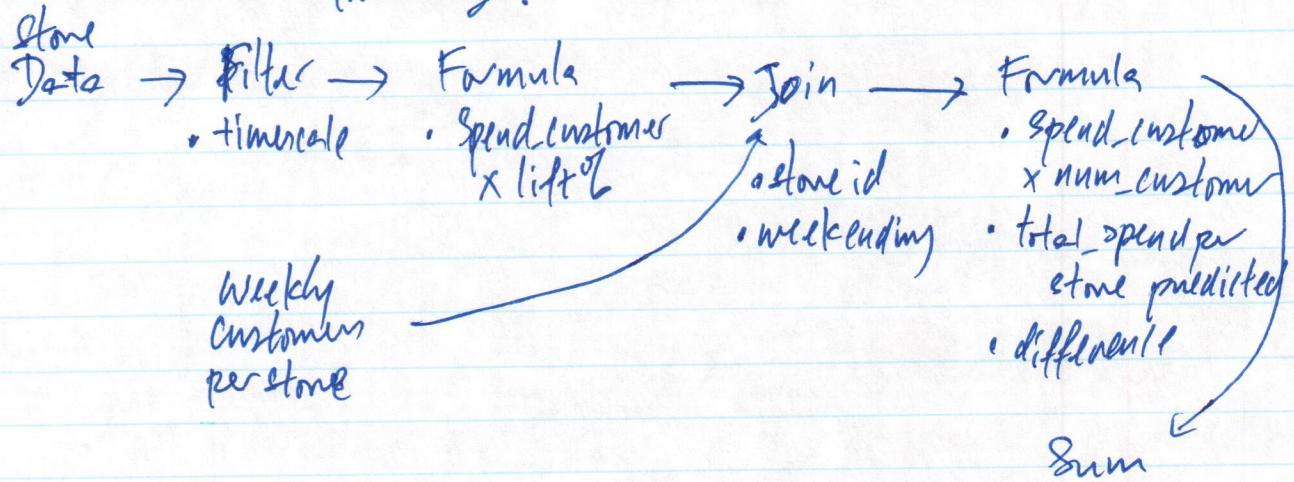
- ~~baseline~~ serve as baseline
- match to treatment unit based on control variables
 - ↳ sales volume
 - ↳ products sold
 - ↳ location - e.g. state, city
- AB Controls tool
 - ↳ KD tree
 - metric - distance
- balance number of control units and average distances



- mean + variance of differences between each control + treatment pair differ from 0



Estimated New Sales from Change:



A/B Test Flow to Calculate Lift + Significance for Store Sales.

Data cleanup:

Cleanup Data and
match up stores

Calculate average sales
for comparable period

Calculate growth for each week for each
store against average sales in
comparable period

Calculate growth:

Match up stores:

Create store list
for significance testing

Average growth by testing
and comparable period

Match store list with
average growth by period

Difference:

Difference of growth between
test and comparable period for each
store

Calculate
lift +
significance:

Test significance of
growth difference
using Welch's t-test

Calculate lift +
lift impact of growth
difference between
each treatment store +
corresponding control store

Match Pair Practice:

Price elasticity of demand

- responsiveness of quantity demanded of a goods / service to a change in its price.

A/B testing analyzes whether a change applied to an experimental group is statistically different from a control group.

Control - epo-facial - \$ 98-99

Treatment - $-\$76.99$) treatment group
 $\$87.99$

control variables

- average sales per store per month
 - average inventory per store per month
 - average number of facets by SKU per store per month
 - average sales of " " " " "

Continuous trends

- Growth in store traffic
 - \ visitors per week
 - Seasonal patterns in sales volume
 - \ total sales per store per week

A/B Trend tool

B trend tool
needs 1 yr of data + minimum 12 periods in addition to year if measuring weekly.

e.g. 1 customer every 10 weeks

→ length of experiment = 10 weeks

→ length of experiment = 10 weeks
→ comparison period = 10 weeks, same as testing period

→ length historical data = 1 yr + 12 periods

$$52 \text{ weeks} + 12 \text{ weeks} = 64 \text{ weeks}$$

Test period, 10 weeks

aggregate weekly sales, invoices

→ weekly per store invoices/traffic \Rightarrow A/B Trend tool

→ Store list data matched of control to treatment stores \Rightarrow A/B Controls

→ store sales, gross margin per week per store \Rightarrow A/B Analysis tool.

weekly store traffic \rightarrow A/B

Trend
Tool

Jain
Tool

Store list
• randomly
assign treatment/
control

Date
Filter \rightarrow Filter control \rightarrow A/B Controls

• Region = Midwest
• Test group = CC

Filter \rightarrow Filter control \rightarrow A/B Controls

• Region = West
• Test group = CC

Filter \rightarrow Filter control \rightarrow A/B Controls

• Region = East
• Test group = CC

Bonus \leftarrow A/B Analysis

control

Filter
Treatment, Test group
= 87.99

Union

Jain

control -
treatment
pairs

Stone sales
• weekly sales per
store

Bonus \leftarrow A/B

Analysis

Filter
• Test group
= 79.99

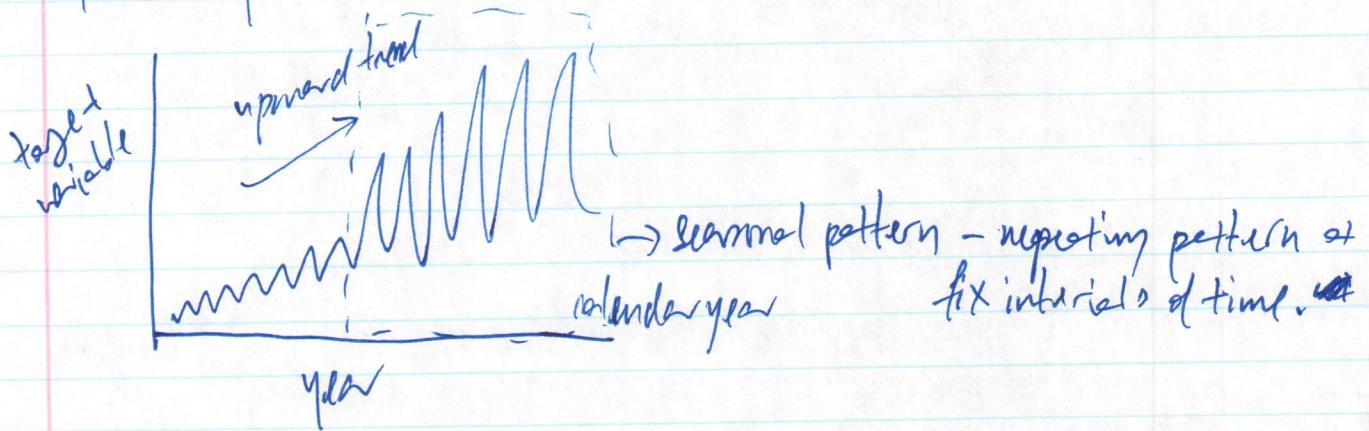
Time Series Forecasting

- Order matters
- continuous time interval
- equal spacing between 2 consecutive measurements

) time series attribute

Steps

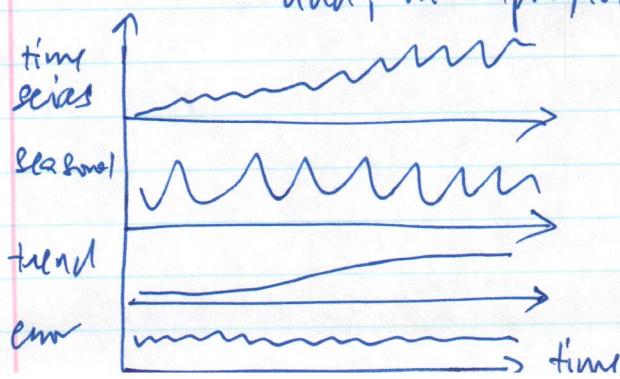
- investigate + clean data
- determine trend, seasonal components
- apply fitted to ARIMA model
- forecast future



- cyclical pattern
- upward/downward trend
 - longer than seasonal
 - much harder to predict
 - e.g. stock market

Exponential Smoothing Models

time series - error, trend, seasonal, add, multiplicative, left out

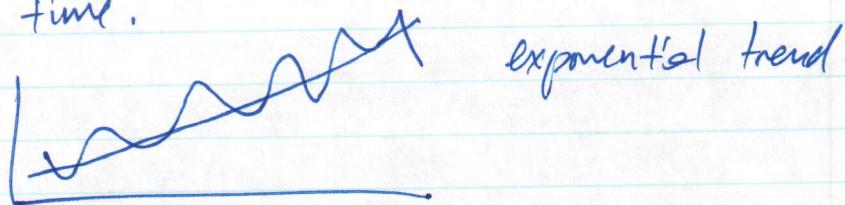


Decompose time series data in seasonal, trend and remainder error components.

Additive - trend, seasonal variations are relatively constant over time.



Multiplicative - trend, seasonal variation increases / decreases over time.



Scenarios

No trend, no seasonal

" , seasonal constant

" , seasonal increasing

Trend-linear, no season

" , seasonal constant

" , seasonal increasing

Trend-exponential, no seasonal

" , seasonal constant

" , seasonal increasing

ETS Models (Err, Trend, Season)

1. Simple Exponential Smoothing

- no trend, no seasonality

$$\text{Forecast} = w_1 x_1 + w_2 x_2 + \dots + w_t x_t$$

$$w = \alpha (1-\alpha)^t$$

α = range 0 to 1

- more weight placed on more recent data

2. Double Exponential Smoothing / Holt's Linear Trend Method

- add trend + level

linear

3. Exponential Trend Method

- multiplicative trend + level
 (exponentially)

4. Damped Trend Methods

- dampens trend line into flat line into future
- additive or multiplicative
- damping parameter, Φ

5. Holt-Winters Seasonal Method

- trend, level, seasonal
- add, multiplicative method
- use with Φ

ARIMA Models:

- Auto Regressive Integrated Moving Average

1. Non-Seasonal ARIMA (p, d, q)

AR - p - periods to lag for in prediction

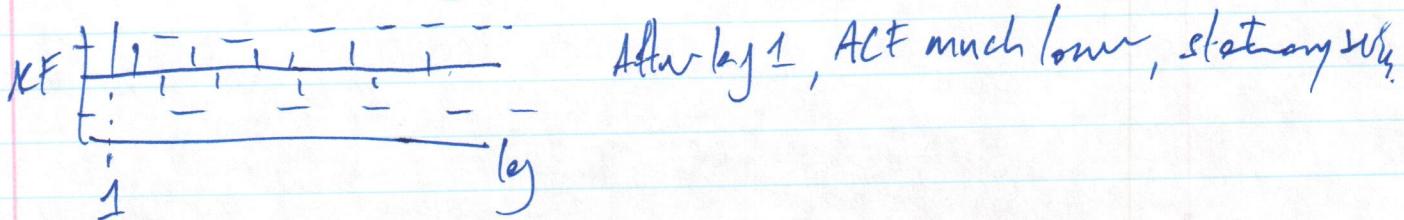
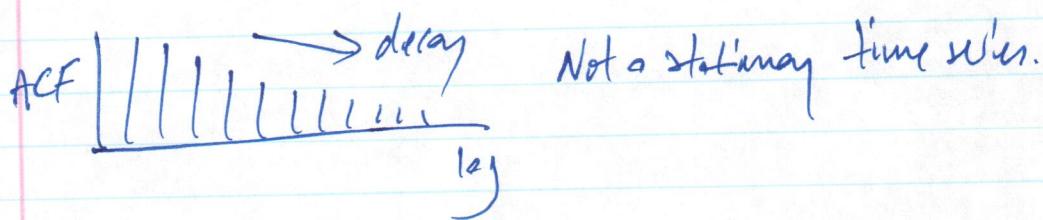
I - d - number of transforms used to transform a time series into a stationary one.

MA - q - lags of error component (not explained by trend/seasonality)

Stationary time series - mean + variance constant over time.

Auto Correlation

- how correlated a time series is with past values.
- Autocorrelation Function (ACF) plot

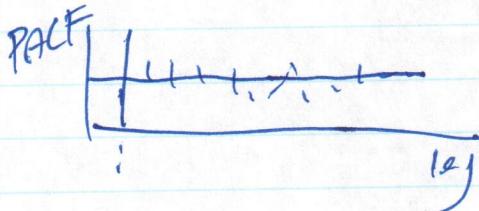


AR model - positive ACF at lag 1.

MA model - negative ACF at lag 1.

- cuts off sharply after a few lags.

Partial correlation - correlation between 2 variables controlling for values of another set of variables.



AR model - PACF drops off after lag-k indicates AR-1c model

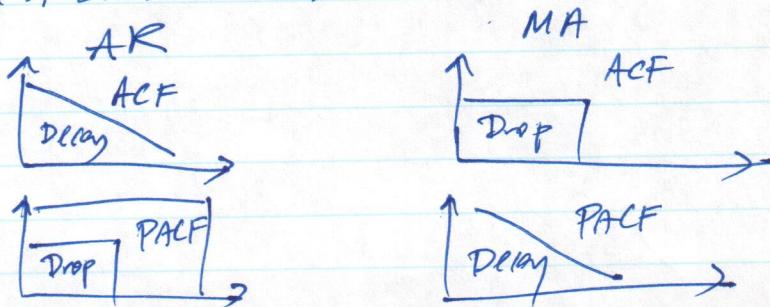
MA model - PACF drops off gradually, indicates MA model

2. Seasonal

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

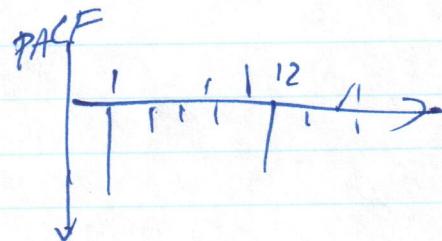
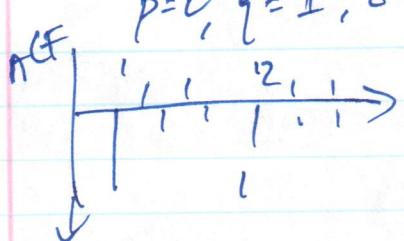
- use seasonal lag to calculate differenced series.

- rest is similar to non-seasonal



$$P=0, Q=1, D=1$$

$$P=0, q=1, d=1$$



P, d, q - look at lags not same as seasonal lags - 12, 24

P, D, Q - only look at lag at seasonal lags - 12, 24

Interpreting Results

Hold out sample - subset of data withheld to check accuracy of prediction of model.

Residuals = difference between observed + forecasted value.

Good models :

- Residuals in model uncorrelated

↳ plot in ACF

↳ adjust for lags

- Residuals means = 0

↳ otherwise subtract mean from model!

Measures of Error

1. Mean error (ME)

↳ average of difference between actual + forecasted.

2. Mean percentage error (MPE)

↳ average of percentage difference between actual + forecasted.

3. Root mean squared error (RMSE)

↳ sample standard deviation of differences between predicted and observed values.

4. Mean absolute error (MAE)

↳ average of absolute errors. less sensitive than RMSE to large errors.

5. Mean absolute percentage error (MAPE)
' compare between 2 different data series.

6. Mean absolute scaled error (MASE)
$$= \frac{\text{mean absolute error of model}}{\text{mean absolute value of 1st difference of series}}$$

' relative reduction in error compared to naive model.

Akaike Information Criterion (AIC)

- measure of relative quality of statistical model
- balance of goodness of fit of model and complexity of model.
- lowest is best

Confidence Interval

- interval within which forecast will lie.

Segment and Clustering

Standardization - treat all units the same through generalization.

- large errors, but carry.

Localizat - treat local needs individually

- small errors, but effort needed to differentiate.

Clustering - statistical methodology to group similar objects into clusters based on multidimensional features or attributes.

Distancie - to quantify or measure similarity.

- Data Preparation

- 1. Selecting data

Retail example for required data

Benchmarking system:

- square footage
- business hours
- management
- street traffic
- employee + hours
- population
- competition

Pricing rules:

- income
- life stage
- distance to competition
- type of competition
- regulatory zones

Product Allocation + Space Planning:

- area demographics - income, age, ethnicity
- product category sales history
- sell-through - compare amount of inventory retailer receives from manufacturer against what is sold to customer.
- product discount trends

Predetermined bias

Retail - e.g. product category sales

- fair categorization alone in segment of similar stores, despite demographic + social data.
- cluster based on demography + historical sales to see any bias present.

- 2. Data types - ordinal
 - hot encoding

3. Data quality - mining data

4. Scale data

- Z-score - number of standard deviations from mean

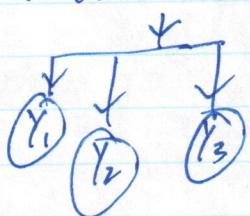
$$Z = \frac{x - \mu}{\sigma}$$

- Standard - rescale to 0 to 1

• Variable Reduction:

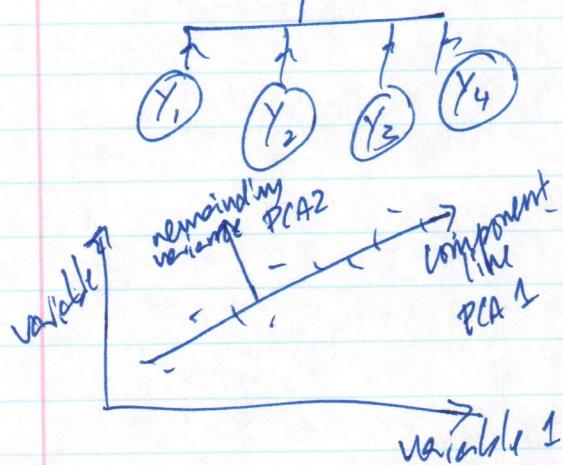
- account for most of variance in all observed variables.

Factor analysis - hidden factor underlying variables in question.

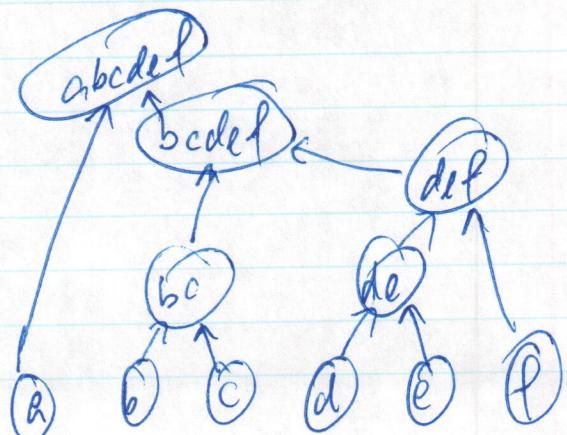


- correlation between variables
- find cause between variables.

Principal Component Analysis (PCA) - summarize all variance within total variables into fewer components, reducing # variables.



- account for total variance explain as much of variance as possible.
- PCA tool



• Clustering Models:

- Hierarchical Clustering
1. bottom-up

2. step down

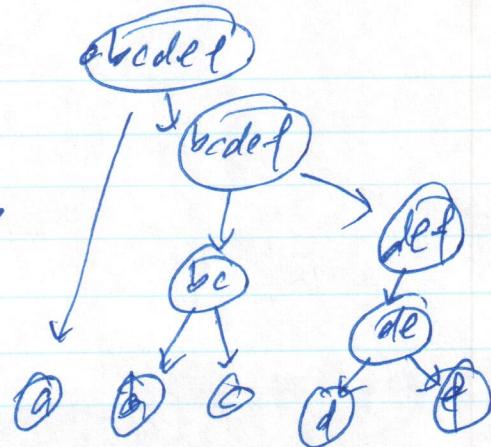
cluster based on linkage

1. Single - ^{closest} point to point distance

2. Complete - farthest point to point distance

3. Average - average distance

4. Centroid - centroid to centroid



Stops when all objects are in its cluster - no need to specify K , # cluster

- K-centroid Clustering

need to specify k

initial centroid chosen randomly

repeat
 → assign object closest to centroid
 → calculate centroid
 until no or little change in centroid

worse better with large number objects
even number of clusters
different density of clusters.

Highly sensitive to outlier, affecting centroid

Validation:

- internal validation

- external validation

evaluate results measure of solution
compensate results using an index.

objective of clustering

1. minimized intra-cluster separation → compactness

2. maximized inter-cluster separation → separation or distinction

Indices

1. Adjusted Rand Index

- how similar objects within a clusters are or cluster stability

2. Calinski - Harabasz Index (CH)

- measures both compactness + distinctiveness of clusters.

K-Centroids

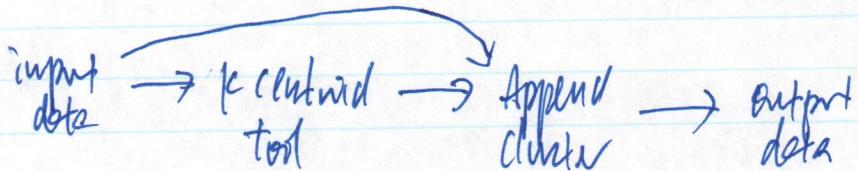
Input → Diagnosis → Browse

Tool

- standardize < z-score
- cluster methods
 - k-means
 - k-medians
 - Neural Gas
- min/max number of clusters

} calculates Adjusted Rand and CH Indices between min and max ranges of # clusters.
AR index - higher, more stability
CH index - higher, more compactness + more distinctiveness.
→ want high median and interquartile range is small.

- vary # clusters, methods, standardize, index metrics.



Validating Models

- Internal validation

 └ best number of clusters

- External validation

 └ hold out set

 └ validate through visualization

Date Visualization in Tableau:

1. worksheet - 1 visualization
2. dashboard - multiple visualizations
3. stories - combination of 2