

EDA

What kind of variables

What distributions

Anomalies?

Interesting associations?

Central location

Mean

Median (artiles x)

Mode (categ)

Quant

Numerical,
Ratio
Interval

Categorical

Nominal

Ordinal

Logical order

Locations

↳ 25th, 50th, 75th percentile

Spread/Variation

↳ diff btw Min & Max Range

Third & first quartile IQR

Variance & s.d (dispersion of data around mean)

Cat or quant
reference values
or dummy
Optimal values

Symm/normal
more than 3 S.D
mean

Shape

uni / bi / complicated model

left skew / right skew (long tail)

Kurtosis (thin or thick tails normal)

Box plot

1.5 IQR 3Q

& 1 less than

1.5 IQR 1Q

Normal dist assumptions → symmetric, defined mean & sd,
bellshaped, $m = md = med$

Distribution → hist, normprob, Q-Q plots,

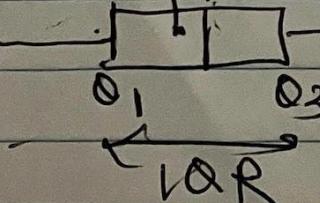
$Q_1 - 1.5 \text{ IQR}$ Mean
Med

Box plots

$Q_3 + 1.5 \text{ IQR}$

asymptotic to x
axis
(bounds are $\pm \infty$)

outliers



outliers

Deviation
from normality

to show range

small s.d. data close to mean
large s.d. → data away from mean or spread out

Why EDA → make necessary assumptions about analysing it properly
→ understanding of variables, relations
→ which stat analysis might be appropriate

Nominal - no logical ordering

↳ while modelling converted dummy variables

↓
indicate presence/absence

Linearize → reference / onehot effect

ML models → one hot / other encoding

↳ dummy variable for each level

→ similar to one hot but reference level is dropped

Interval variables → arithmetic operations ✓

Continuous → quantitative can take any value within range

Categorical → qualitative, represent types or categories

ordinal variables → treated as interval variables / categorical variables

histogram - continuous

barplot → nominal

$$S.D = \sqrt{\text{variance}}$$

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \text{mean})^2}$$

covida

Mean - avg

Median - middle

mode - most common

$$IQR = 25^{\text{th}} \text{ perc} - 75^{\text{th}} \text{ perc}$$

→ S.D preferred same units
of variance ✓

Anomalies → normal dist 3.S.D from mean

Boxplot → $> 1.5 \times IQR$ Q3
 $< 1.5 \times IQR$ Q1

Point Estimates sample statistic \rightarrow estimate of pop parameters
so... to know more about pop parameters

Variability among Samples

CLT \rightarrow dist of sample mean

is normal regardless of pop dist's shape, if sample size is large enough.

standard error \rightarrow variability from estimate (S_x)

standard deviation \rightarrow variability in data

$$S_x = \frac{S}{\sqrt{n}}$$

provide info variability around statistic

CI \rightarrow contains pop parameter with some degree of confidence

CL \rightarrow % indicates strength of confidence that interval created actually captures true parameter

95% created intervals capture pop parameter

exp & population
Hypothesis \rightarrow if true pop value sig diff from assumed value

statement we wish to investigate scientifically

through process of stat inference

- ① initial assumption, Null to
- ② analyze data to determine whether observations are likely given null hypo

- ③ if given data observed data unlikely reject initial assumption favour of alternative

$\alpha \rightarrow$ sig level
how much evidence to reject null hypo

$p \rightarrow$ probability that we observed something as extreme or more extreme than we did

$p < \alpha$ reject null

Type I \rightarrow reject null

Type II \rightarrow fail to reject null

t test

mean

One sample t test

mean against hypo value

Two side t test

t can be +ve or -ve

hypo \rightarrow inside CI

fail to reject null

outside CI

reject null

One side test

sale price $> 170K\$$

shape, var, symmetric

* Stat significance \rightarrow observed data is not due to random chance

Two sample t test

testing diff btw 2 group means

assump ind observations

normality each group

equal var

QQ, Shapiro

F test (normality)

2 sample test \rightarrow with equal
without equal

When normality fails

Wilcoxon rank sum +

Man with they wilcoxon test
diff in means/medians/dist

domain

F test, Shapiro, t test

null \rightarrow equal mean
equal var
normality

Predictor	Catg	Cont	cont/Catg
Cont	ANOVA	OLS	OLS
Catg	Log	Log	Log

Honest Modelless *

Linear Model

↳ explanatory

Response

→ explore dist, outliers, missng

→ then split before relationship

↓
train / valid / test

Overfit *

↳ model captures noise on train

→ pattern not hold in validation or test

→ performance suffers

doesn't generalize
model overfit

pred

→ x_i can predict y

→ predictions focus

→ many variable &
complex

50 - 40 - 10 lots

70 - 20 - 10 not enough
cross

Bivariate EDA

expected value of one variable
changes at diff levels of another

② grouped box plots

overhead hist

cont vs category

① linear association scatter plot
straight line

One way Anova

$H_0 \rightarrow$ means of each level equal

categ input & quant response $H_A \rightarrow$ at least one diff

Null \rightarrow means equal

All \rightarrow means not equal

\Rightarrow we do reference level coding

$$Y = \beta_0 + \beta_A X_A + \beta_B X_B + \epsilon$$

$$C_i \rightarrow Y = \beta_0$$

β_A diff means for A vs for C

(1) EDA

(2) verify Ass

(3) P value of overall F test in Anova

$P < \alpha$ reject null

Assump \rightarrow Ind observation \rightarrow good data collect

\rightarrow groups normally (residuals) \rightarrow QQ, hist / formal test

when normality fails Kruskal Wall's \rightarrow equal var \rightarrow classic \rightarrow formal test (homoskedasticity)

\rightarrow Welch \rightarrow Leven's \rightarrow fligner's

require norm nonnormality

Posthoc Anova

\hookrightarrow which groups diff, multiple sample t-test

\downarrow multiple hypothesis

Ex post Error Rate

Comparisons each having error rate

prob that I make atleast one error

$$1 - (1 - \alpha)^n$$

\downarrow compound error

to control \rightarrow Tukey (all pairwise)

Punnett (comparisons to control group)

cont resp & continuous pred.

corr(-vetodre)

SLR

→ Corrcoeff → measures
Univrelation btw 2
continuous attributes

→ if we don't explore data
sometimes anomalous observations
declare, relationships are
significant → outliers skew

→ raw correlation misleading
so scatterplot

corr value, test, plot

→ can mislead when both
variables affected by other variable.

→ how much one variable
changes as other changes by
one unit

→ creates Lineareq by
Minsquares btw
observed data & model
predictions

OLS → Min SS residuals
Cost function

→ predict expected value of
y for each value of x

→ y is expected to
change for unit chg

purpose

$$y = \frac{\beta_0 + \beta_1 x + \epsilon}{\text{determin} \downarrow \text{random}}$$

$$H_0: \beta_1 = 0 \quad \beta_1 \neq 0$$

global Ftest

param + Hot

pearson

equal

Assump

linearity of mean

errors are normal

no pattern
in residual plot

Hist QQ, normally

ResVspxd nicebox

data collection
& explore for

auto-correlation

equalVar

Independent

2 way

n way Anova

2 ways → 2 variables

K₁ & K₂
categories

Model → relation b/w exp var &
Res var

Effect → expected chg in res var
w/ value f explanatory

2 ways
1) EDA
2) 2 way anova → each variable
has F test

post hoc test

3) Main effect → effect f single
explanatory var

Interactions → effect f one variable
changes as levels f another

mask effects of variables

~~slicing~~ → f test mean f one
within effects testing }

ANOVA
essentially
a linear

assump → Indep observations

→ Eq of variance (various)

→ Norm of categories (various)

Multiple Linear Reg

linear
comb of variables

Multicollinearity

→ more than one predictor

$\beta \rightarrow$ predicted change in Y

with one unit increase in X
given all other variables
held constant *

→ pred var are

correlated with
each other

→ Linear reg breaks only
when collinear perfect

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$$

$$\beta_0 + \beta_1 x_1^2 + \beta_2 x_2^2 + \epsilon$$

$$\beta_1 = \beta_2 = \dots = 0$$

Global → None of var useful

F test → at least one variable
useful

If Global F test sig →
Local test

Adv

investigate
relation among Y &
several indp variables

Disadv

↑ complexity makes
difficult

→ which model
best
interpret models

Assump Linear Reg

(1)

mean of Y accounted
modeled by linear fun of X

(2)

random error term assumed
to have normal dist mean 0

(3)

const var

(4) errors independent

Application
↓
Predict

$R^2 \rightarrow$ prop of variance in dependent var explained by independent var

focus on predictions

value of coeff & stats less important

$R_{adj}^2 \rightarrow$ consider no variables & permutes

excessive variables

→ Explain to develop underst
pred Var & ResVar

* with addition of variables $R^2 \uparrow$ in MLR so $R_{adj}^2 \rightarrow$ penalizes model for adding variable that doesn't provide any useful information

Categorical variables in MLR → dummy encoded

Avg diff between catg Y & N α $1,0 = X$
 $B_0 + \beta_1 X$

effect coding
average diff between catg Y & overall avg cat Y & N.

Model Selection

Information \rightarrow used to select variables for model

Selection Algorithm \rightarrow auto techq to evaluate variables
on some selection criteria

↓
Stepwise \hookrightarrow All or nothing selection (R^2 , R_{adj}^2 , Mallows(p))
From back, stepwise

Model selection based on training data (to avoid overfitting)

AIC \rightarrow crude, large sample approx f. known one out CV

BIC \rightarrow favors small models / penalizes model complexity more.

\rightarrow low value better than higher, no amount of lowering better.

forward start with intercept
build the model one variable at a time

} backward
start with full model
& remove variables

dir=forward direction both

stepwise
start with empty
& add & remove

K = 2 AIC

$K = \log(n) + 2k$ (train & dev)

BIC selection

$K = -2 \ln L(\theta) + k \ln n$ (train & dev)
 \hookrightarrow p-values selection with alpha 0.05

Issue with AutoSel Algo

- bias in parameter estimates, pvalue, standard errors
- Incorrect cal & DDF
- pvalue tend to err on side of overestimating significance.
(Type I error)
- local best not global
- don't use result from ASel algorithm blindly

conservative pvalues

sample size ↑
pvalue ↓

→ Automated stepwise can provide subset of potential variables

→ no model blindly finished

explore other models &
investigate assumption

→ Adjust pvalues for large samples.

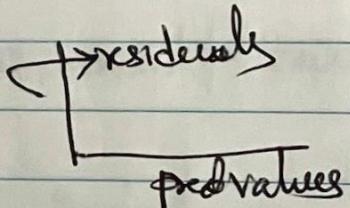
Diagonistics

1. Reg assumptions

Violated

- (1) mean of Y accurately modeled by linear function X → misspecified model not monotonic
- (2) constant var random error → S.E. compromised
coeff parameter estimates ✓
- (3) random error \rightarrow normal dist with mean zero → parameter estimates ✓ effects test results
- (4) errors are independent → no effect parameter estimates standard error compound
- (5) no perfect collinearity

Residual plot \rightarrow randomly scattered



Misspecified \rightarrow pattern detected
(poly, int, "spline", ...)

* Model hierarchy when adding higher order, all lower terms included in model

when straight line not appro

- \rightarrow fit poly / complex reg model
- \rightarrow transform dep / indep var to obtain linearity
- \rightarrow fit non linear reg model,
- \rightarrow fit non param reg model (LOESS)

Const variance (homoscedacity)

→ F test

→ examine residual plots or spearman rank correlation

→ H₀ violated (Hypothesis, E, CI)

but don't t/f, χ^2 not valid

0 → Variance is not homo

corr coeff between

+ve $\sqrt{\uparrow}$ mean \uparrow

absolute value of residuals

-ve $\sqrt{\downarrow}$ mean \uparrow

& predicted values

H₀ relation is not linear test will not detect it

→ accounting hetero

[use WLS]

transform data

use diff distribution

Normality

→ histogram of residuals

→ Normal probplot (Q-Q plot)

→ formal tests for normality

accounting for lack of normality

→ scatterless (robust reg)

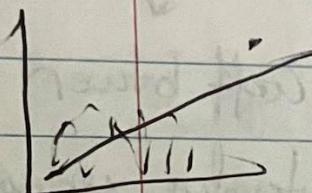
→ nonnormal (transf needed)

can by BoxCox transformation

Autlier → large standard residuals, far away from Y direction

Influential residuals, far away from X direction

Leverage points → points fall outside normal range X space
large influence on regression line.



Residual only focus on outliers

Diagnostics → standard studentized { good for outliers
COOK'S D DFFITS DFBETAS } good for influential observations
hat values

how to handle Influential observations

→ Recheck data to ensure no transcription or data entry

→ data is valid model adequate

↓
run analysis with & without influential observations
high order terms poly & interaction
nonlinear model

→ Robust Regression

→ WLS

Collinearity

- look at corr matrix
- use VIF (> 10 potential coll
 > 5 GVIF...)
- dealing
 - exclude red independent variables
 - Redefine variables
 - use biased Reg techniques
 - center independent var in poly reg models
or model with interaction terms.

Model building & Scoring for prediction

Linear reg BLUE

best linear unbiased estimate

→ If assumptions hold best linear min variance of all unbiased estimators.

If assumptions don't hold?

If biased estimators had small variance?

no. of var ↑ → assumption

multicollinearity

coeff vary

leads to overfitting

high variance

→ More var than observations

Regular Reg puts constraint on estimated coeff & shrink estimates to 0

coeff biased but improve variance

→ Reg Reg add penalty term to $\text{min}(\text{SSE} + \text{penalty})$

loss function

$\lambda = 0$ OLS

$\frac{1}{2}$ penalty alpha=0 ridge → approach to 0

$\lambda \propto \text{coeff}'$

L_1 penalty

alpha=1

Lasso → coeff = 0

variables removed

Lasso → does variable selection

Ridge → keeps all variables

ElasticNet → Lasso + Ridge alpha between 0 & 1

at

Fair of Overfitting → need to select λ for any regularised regression approach

→ don't want to minimize variance to point of overfitting to training

Cross-validation prevent overfitting

<u>Model metrics</u>	RMSE	not easily interpretable
	MAE	not scale invariant
	MAPE	not symmetric

Categorical Data Analysis

- detect freq. of data
- possible association among variables

when sample size
small or differ
can't compare

crosstable function in R
chisquaretest

Nell

Pearson

any type of
likelihood

Azhaaratu

Wilkelihood

categorical
variable.

Sample size req

80% or more

expected
cellcount > 5

when we don't meet assumptions

↳ fisher's exact test

ordinal more info → Mantel Haenszel chisquare test

strength of association

[odds ratio, Cramér's V

odds = $\frac{P}{1-P}$

Spearman's correlation]

how much more likely w.r.t odds a certain event occurs in one group relative to its occurrence in another

oddsratio

Cramers V more than 2 categories in one or both variables use Cramers V

odd's ratio B vs B

Spearman's strength of association btw 2 ordinal variables

→ Pearson on rank instead of value of observations

Continuous Target → Linear Reg

Anova

Log Reg Reg

Categorical Target → Logistic Reg

Intro to Logistic Reg

↳ predict probability between 0 & 1

→ parameter estimates don't enter model equation linearly

→ rate of change of probability w.r.t. X varies

↳ logit applied to probability $\log\left(\frac{p_i}{1-p_i}\right)$

↳ relation b/w logits linear $\log \text{ applied to probability}$

→ logits unbounded

assumption $\begin{cases} \rightarrow \text{Indp observations} \\ \rightarrow \text{logit is linearly related} \\ \quad \text{to variables} \end{cases}$

→ $100 \times (e^\beta - 1) / \beta$. change in odds

→ β is change in logit

~~We want high concordant~~ Concordant $(0, 1)$ \rightarrow high pred prob

Discordant $(0, 1)$ \rightarrow high pred prob

~~Concordant~~ Tied $(0, 1)$ bonus eligible has same pred prob as non bonus eligible

($\text{glm}(\text{ } \text{ })$ function)

Variable selection

↳ for, back, stepwise, Lasso

→ Reg Reg can use same link function to obtain logistic
Regression