

Contents

1	Basic Terminologies	6
1.1	Population	6
1.2	Sample	6
1.3	Histogram	6
2	Law of large numbers	7
2.0.1	Weak Law of large numbers	7
2.0.2	Strong Law of large numbers	7
3	Central-Limit Theorem	8
4	Measure of Central Tendency	8
4.1	Mean/Expected Value	8
4.2	Median	8
4.3	Mode	8
5	Measure of Spread	8
5.1	Variance	8
5.2	Standard Deviation	9
5.3	Inter Quartile Range(IQR)	9
6	Probability Distributions	9
6.1	Discrete Distributions	9
6.2	Continuous Distributions	10
7	Discrete Probability Distributions	10
7.1	Benoulli Distribution	10
7.1.1	Probability Mass Function(PMF)	10
7.1.2	Statistical Parameters	10
7.1.3	Examples	11
7.2	Binomial	12
7.2.1	Probability Mass Function(PMF)	12
7.2.2	Statistical Parameters	12
7.2.3	Examples	13
7.3	Negative-Binomial	13
7.4	Multinomial	13
7.5	Geometric	13

7.6	Hypergeometric	13
7.7	Poisson	13
7.8	Discrete Uniform	13
7.9	Normal Distribution	13
7.9.1	Use cases	13
8	Shape fo the Distributions	13
8.1	Skewness	13
8.1.1	Right-skewed	13
8.1.2	Left-skewed	13
9	Hypothesis Testing	14
9.1	Null Hypothesis(H_0)	14
9.2	Alternate Hypothesis(H_a)	14
9.3	One sided Test	14
9.4	Two sided Test	14
9.5	Test Statistic	14
9.6	Sampling distribution under H_0	14
9.7	p-value	15
10	Machine Learning	15
10.1	Type of Machine learning	15
11	Types of Data	15
11.1	Qualitative Data	15
11.1.1	Nominal	15
11.1.2	Ordinal	16
11.2	Quantitative Data	16
11.2.1	Discrete	16
11.2.2	Continuous	16
12	Data Clearning	16
12.1	Define the target variable	16
12.2	Policy Document	16
12.3	Outlier Detection	17
12.3.1	Inter Quartile Range(IQR)	17
12.3.2	Z-score	17
12.4	Percentile/Quantile Trimming	17

12.5	Domain Rules	18
12.6	Impute missing values	18
12.6.1	Numerical data	18
12.6.2	Categorical	18
12.6.3	Time series	18
12.6.4	Images	18
12.6.5	Text	19
12.7	Data transformation	19
12.7.1	Standardization	19
12.7.2	Min-Max scaling	19
12.7.3	Robust scaling	19
12.8	Log transformation	20
12.9	Text processing	20
12.9.1	Lowercasing	20
12.9.2	Noise Removal	20
12.9.3	Tokenization	20
12.9.4	Stemming	20
12.9.5	Lematization	21
12.10	Handling spelling and special entities	21
12.11	Name Entity Recognition	21
12.12	Encoding	21
12.12.1	Bag of Words(BOW)	21
12.12.2	Term frequency-Inverse document frequency (TF-IDF)	21
12.12.3	Embedding	22
12.13	Image processing	22
12.13.1	Resizing & Cropping	22
12.13.2	Normalization/Scaling	22
12.13.3	Grayscale conversion	23
12.13.4	Data Augmentation	23
12.13.5	Noise Recution/Smoothing	23
12.13.6	Histogram Equalization	23
12.13.7	Color Standardization/White balance	23
12.13.8	Segmentation	23
13	Feature engineering & selection	23
13.1	Filter Methods	23
13.1.1	Correlation coefficient	24
13.1.2	Chi square test	24

13.1.3	ANOVA F-test	24
13.1.4	Variance Threshold	25
13.2	Wrapper Methods	25
13.2.1	Forward selection	25
13.2.2	Backward Elimination	25
13.2.3	Recursive Feature Elimination(RFE)	26
13.3	Embedded Methods	26
13.4	Dimensionality Reduction	26
13.4.1	Principal Component Analysis (PCA)	26
13.5	Domain Knowledge	26
14	Supervised Learning	26
15	Unsupervised Learning	26
16	Time Series	26
16.1	Properties of Time series	26
16.1.1	Order	26
16.1.2	Time dependence	27
16.1.3	Stationarity	27
16.1.4	Trend	27
16.1.5	Seasonality	27
16.1.6	Cyclic	27
16.2	Suitability Test for TSA	28
16.2.1	Time dependent	28
16.2.2	Auto Correlation	28
16.2.3	Stationary	28
16.2.4	Data Quality	29
16.2.5	White noise	29
17	Tests for time series	29
17.1	Augumented Dicky Fuller Test (ADF)	29
17.2	Unit Root	30
17.3	Auto Correlation Function(ACF)	30
17.4	PACF	31
17.5	Time series decomposition	31
17.6	Periodogram	31
17.7	ARCH test	31

18 Evaluation Metrics for Regression	31
18.1 Mean Absolute Error	31
18.1.1 Gradient of MAE	31
18.2 Mean Squared Error	32
18.2.1 Gradient of MSE	32
18.3 Root Mean Squared Error(RMSE)	32
18.3.1 Gradient of RMSE	32
18.4 R-Square(R^2)	33
18.5 Adjusted R^2	33
18.6 Mean Absolute Percentage Error(MAPE)	34
18.7 Huber Loss	34
18.7.1 Gradient of Huberloss	35
19 Evaluation Metrics for Classification	35
19.1 Basic Terminologies	35
19.1.1 True Positive	35
19.1.2 False Positive	35
19.1.3 True Negative	35
19.1.4 False Negative	35
19.2 Accuracy	36
19.3 Precision	36
19.4 Negative Predicted Value (NPV)	36
19.5 Recall(TPR)	37
19.6 Specificity	37
19.7 False Positive Rate(FPR)	37
19.8 F-score	37
19.8.1 Derivation	38
20 AUC-ROC curve	38
20.1 Receiver Operating Characteristic (ROC)	38
20.2 Area Under the Curve (AUC)	38
21 Regularization	39
21.1 L1	39
21.2 L2	39
21.3 Elastic net	39
22 TODOO	39

Statistics For Data Science

Akash Tesla

July 2025

1 Basic Terminologies

1.1 Population

An entire set of items you want to study

1.2 Sample

A subset of population used to estimate statistical behavior of the whole population

1.3 Histogram

A histogram is a graphical representation of numerical data that groups the data into bins and displays the frequency of data points within each bin as bars

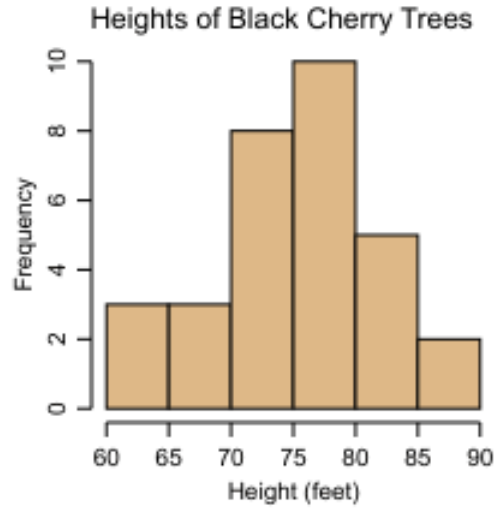


Figure 1: Example of a Histogram

2 Law of large numbers

As the number of trials (or samples) increases, the sample average (or empirical mean) will converge to the expected value (or population mean).

2.0.1 Weak Law of large numbers

The weak law states that the sample average of a sequence of independent identically distributed (i.i.d.) random variables converges in probability to the expected value as the number of samples goes to infinity

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu \quad \text{as } n \rightarrow \infty$$

which means,

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} \mathbf{P}(|\bar{X}_n - \mu| > \varepsilon) = 0$$

2.0.2 Strong Law of large numbers

The strong law states that the sample average of a sequence of i.i.d. random variables converges almost surely to the expected value as the number of samples goes to infinity

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mu \quad \text{as } n \rightarrow \infty$$

Which means,

$$P(\lim_{n \rightarrow \infty} \bar{X}_n = \mu) = 1$$

3 Central-Limit Theorem

4 Measure of Central Tendency

4.1 Mean/Expected Value

Average of all data points, sensitive to outliers since a single large outlier could easily skew mean

$$\mu = \frac{\sum x_i}{n}$$

4.2 Median

The middle data point when data are sorted, robust to outliers

4.3 Mode

The most frequent data point of the dataset

5 Measure of Spread

Range: Difference between minimum value and maximum value

$$Range = x_{max} - x_{min}$$

5.1 Variance

Average squared deviation, Variance represents Expected variance between mean and data points, It's basically MSE of a model that just predicts mean, that kinda gives an intuitive understanding of how it measures spread

$$\begin{aligned}
\sigma^2 &= E[(X - \mu)^2] \\
\sigma^2 &= E[(X - E[X])^2] \\
\sigma^2 &= E[X^2] - (E[X])^2 \\
\sigma^2 &= \frac{\sum (x_i - \mu)^2}{n} (Population) \\
s^2 &= \frac{\sum (\bar{x}_i - \mu)^2}{n - 1} (Sample)
\end{aligned}$$

5.2 Standard Deviation

Root of Variance, RMSE of a model that just predicts mean, standard deviation gives in interpretable terms like RMSE

$$\sigma = \sqrt{\sigma^2}$$

5.3 Inter Quartile Range(IQR)

Difference between 75th Percentile/3rd Quartile and 25th Percentile/1st Quartile, it is used for outlier detection

$$IQR = Q_3 - Q_1$$

We calculate lower bounds and upper bounds to detect the outliers

$$\text{lower bound} = Q_1 - 1.5 \times IQR$$

$$\text{upper bound} = Q_3 + 1.5 \times IQR$$

the data points which values outside of the bounds is considered to be outliers, for more extreme detection $3 \times IQR$ is also used

6 Probability Distributions

6.1 Discrete Distributions

A discrete probability distribution describes the probability of occurrence of each value of a discrete random variable

- Discrete random variable: Countable values like 1,2,3
- Each individual value has an associated probability
- The sum of probabilities for all possible values is 1

$$\sum_i P(X = x_i) = 1$$

6.2 Continuous Distributions

7 Discrete Probability Distributions

7.1 Benoulli Distribution

The benouli distribution is a discrete probability distribution for a random variable which takes only two possibilities, Sucess or a failure

7.1.1 Probability Mass Function(PMF)

$$P(X = x) = \begin{cases} p & \text{if } x=1 \\ 1 - p & \text{if } x=0 \\ 0 & \text{Otherwise} \end{cases}$$

Also written as

$$P(X = x) = p^x(1 - p)^{1-x}, \quad \text{for } x \in \{0, 1\}$$

7.1.2 Statistical Parameters

Mean

Mean is the expected value over many repetitions of the same single-trial experiment, thus it would be p since, p is probability of 1 appearing and (1-p) is probability of 0 appearing

$$\mu = 1 \times (p) + 0 \times (1 - p)$$

$$\mu = p$$

Variance

Variance can be defined as $\sigma^2 = E(X^2) - (E(x))^2$, Refer Variance chapter. For Bernoulli distribution, $E(X^2) = p$, $E(X) = p$, substituting we get

$$\sigma^2 = p - p^2$$

$$\sigma^2 = p(1 - p)$$

Mode

Mode for Bernoulli would be what ever the outcome which is more favored, which can be defined as

$$Mode = \begin{cases} 1 & \text{If } p > 0.5 \\ 0 & \text{If } p < 0.5 \end{cases}$$

7.1.3 Examples

- Will it rain tomorrow?
- Will this patient recover?
- Will this product be defective?

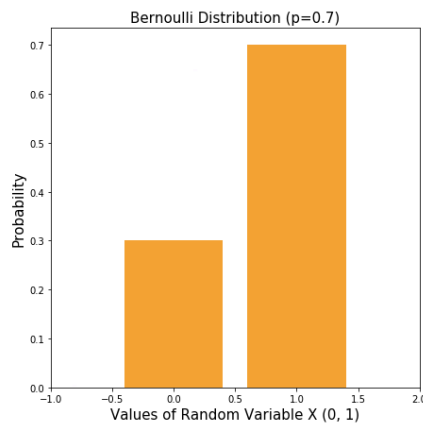


Figure 2: Example of a Bernoulli Distribution

7.2 Binomial

Binomial Distribution is a discrete probability distribution that models the probability of obtaining a specific number of successes in a fixed number of independent trials(n), these independent trials are just Bernoulli trials, you could see the similarity between them in statistical parameters

7.2.1 Probability Mass Function(PMF)

$$P(X = x) = nCx \times p^x \times (1 - p)^{(n-x)}$$

where,

n - no of trials,

p - probability of success

x - number of success

7.2.2 Statistical Parameters

Mean

Mean represents Average number of success from your trials which would be number of trials (n) times probability of success (p)

$$\mu = n \times p$$

Variance

Variance represents Expected variance between mean and data points,

$$\sigma^2 = n \times p \times (1 - p)$$

Mode

$$Mode = \begin{cases} \text{floor}((n+1)p) & \text{if } (n+1)p \text{ is not an Integer} \\ \text{floor}((n+1)p), \text{floor}((n+1)(1-p)) & \text{if } (n+1)p \text{ is an Integer} \end{cases}$$

$$\text{Mode(if } p = 0.5) = \begin{cases} \frac{n}{2} & \text{if } (n+1)p \text{ is not an Integer} \\ \frac{(n-1)}{2}, \frac{(n+1)}{2} & \text{if } (n+1)p \text{ is an Integer} \end{cases}$$

7.2.3 Examples

- How many patients will recover out of 50?
- How many rainy days this month?
- How many defective products in a batch of 1000?

7.3 Negative-Binomial

7.4 Multinomial

7.5 Geometric

7.6 Hypergeometric

7.7 Poisson

7.8 Discrete Uniform

7.9 Normal Distribution

7.9.1 Use cases

- When there is only one trial
- When the outcome is binary True/False Yes/No

8 Shape of the Distributions

8.1 Skewness

Measure of Asymmetry

8.1.1 Right-skewed

tail on the right ($mean > median$)

8.1.2 Left-skewed

tail on the left ($mean < median$)

9 Hypothesis Testing

9.1 Null Hypothesis(H_0)

Null Hypothesis is the default claim basically means no effect/ no difference

9.2 Alternate Hypothesis(H_a)

Alternate Hypothesis is the hypothesis that u want to prove

9.3 One sided Test

When you have to test if the parameter is greater than or less than the hypothesised value, but not both

Null: $H_0 : \mu = \mu_0$

Alternate: $H_a : \mu > \mu_0$ or $\mu < \mu_0$

9.4 Two sided Test

When you have to test if the parameter is different from the hypothesised value in either direction Null: $H_0 : \mu = \mu_0$

Alternate: $H_a : \mu \neq \mu_0$

9.5 Test Statistic

A test statistic is a function of the sample data that is used to decide whether to accept/reject the null hypothesis

$$\text{Test Statistic} = \frac{\text{How surprised we are}}{\text{How surprised we can be}}$$

$$\text{Test Statistic} = \frac{\text{Observed Value} - \text{Expected value under } H_0}{\text{Standard Error of observed value}}$$

9.6 Sampling distribution under H_0

If null hypothesis is true, what would the distribution of my test statistic look like across repeated samples. the sample mean/ test statistic follows normal distribution thanks to CLT (central limit theorem), we use that to calculate p-value.

1. Calculate sample statistic(\bar{x}, s^2)
2. Compute test statistic(t-test,z-test...)
3. Compare the computed test statistic to the corresponding distribution to get the p-value

9.7 p-value

The p-value is the probability of observing your data assuming the null hypothesis is true

if p is small($p < \alpha$), we reject null's hypothesis if p is high($p > \alpha$), we reject alternate hypothesis

10 Machine Learning

Machine learning(ML) is a way of teaching computers to learn patterns from data and make prediction.

10.1 Type of Machine learning

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning
- Semi-Supervised Learning

11 Types of Data

11.1 Qualitative Data

Describes Qualities, Characteristics , or categories

11.1.1 Nominal

Pure categories without order, Example: blood type(A,B,AB,O), brand names

11.1.2 Ordinal

Categories with meaningful order, Examples: Rank, Survey rating

11.2 Quantitative Data

Measureable quantities, Numbers have meaningful terms in terms of magnitude

11.2.1 Discrete

Countable values, no in-betweens. Examples: number of cars

11.2.2 Continuous

Continuous measurements; can take any value within a range, Examples: Height, weight, temperature

12 Data Cleaning

12.1 Define the target variable

Analyse over all data and pinpoint what you want to predict in the dataset,
if your data is discrete like class - classification,
if you have a continuous variable - Regression

12.2 Policy Document

Policy document is a meta data document that defines each column, A policy document could define type of the data, nullable, pattern of the data, Range, units, Logical constraints, description... etc Example:

Column	Type	Null	Pattern/Range	Constraint	Notes
user_id	string	No	regex [A-Z0-9]8	unique	id
age_years	int	No	[0,120]	–	cap outliers
gender	cat	Yes	{M,F,O}	–	harmonize
signup_date	date	No	YYYY-MM-DD	churn_date	UTC
churn	binary	No	{0,1}	target	def: 30d leave
income_inr	float	Yes	[0,1e8] INR	winsorize top 1%	currency check
email	string	Yes	regex email	–	PII → hash

12.3 Outlier Detection

12.3.1 Inter Quartile Range(IQR)

Difference between 75th Percentile/3rd Quartile and 25th Percentile/1st Quartile, it is used for outlier detection

$$IQR = Q_3 - Q_1$$

We calculate lower bounds and upper bounds to detect the outliers

$$\text{lower bound} = Q_1 - 1.5 \times IQR$$

$$\text{upper bound} = Q_3 + 1.5 \times IQR$$

the data points which values outside of the bounds is considered to be outliers, for more extreme detection $3 \times IQR$ is also used

12.3.2 Z-score

How far your point is away from the mean in terms of standard deviation

$$z_i = \frac{x_i - \mu}{\sigma}$$

if $z_i > 3$, the point is very unusual and a potential outlier considers all data above 3 SD as outliers and eliminates it

12.4 Percentile/Quantile Trimming

Trim of top and bottom 1-5 percentage of data, commonly used in competitions to make sure there are no outliers

12.5 Domain Rules

Domain rules like ages should range between 0 - 120, or temperatures should range between -50 to 60 degrees

12.6 Impute missing values

You can either delete the row entirely or chose to impute the missing values

12.6.1 Numerical data

- small missing values - mean
- skewed - median
- important values - regression imputation
- many missing - add a missign flag (0) or remove the column entirely
- if you are not sure you can test all the methods with cross validation and chose a impuation method (recomended for large datasets)

12.6.2 Categorical

- Low cardinatlity(few categories) - Mode impuation
- High cardinatlity - Add a new "missing" category

12.6.3 Time series

- Short gaps - forward/backward fill
- Long gaps - Interpolation
- seasonal data - seasonal average + Iterpolation(optional)

12.6.4 Images

take mean or median of the neighbours to fill in the missing pixels, or drop the image from the dataset

12.6.5 Text

Drop the entire row or add [missing] token

12.7 Data transformation

12.7.1 Standardization

$$x_i = \frac{x_i - \mu}{\sigma}$$

- used in linear regression, SVM, PCA, K-means
- makes mean 0, and std 1
- used to standarize all the numerical featuers so that they don't dominate over one another

12.7.2 Min-Max scaling

$$x_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

- Used in: neural networks, gradient based models, image pixel scaling
- Fits everything in [0,1]
- Helps model converge faster

12.7.3 Robust scaling

$$x_i = \frac{x - \text{median}}{\text{IQR}}$$

- Used in Models sensitive to outliers: regression, SVM, KNN
- Ignores extreme values, centers around median
- Robust to outliers

12.8 Log transformation

$$x_i = \log(x + c)$$

- used in skewed data(income , population, counts)
- used in exponentialisue data to make it more linear
- Compresses large numbers, spreads out small ones to make distribution closer to normal

12.9 Text processing

lemmatizaiton and so on

12.9.1 Lowercasing

Converts all the letters into lowercase

12.9.2 Noise Removal

Remove punctuation, numbers, symbols,stopwords(if not useful)

12.9.3 Tokenization

- Convert sentence into smaller units token
- example: I like data science = [I,like, data, science]
- Tokens are mostly words but not always

12.9.4 Stemming

- chops suffixes from the words
- example: "playing" → "play", "studies" → "studi"

12.9.5 Lemmatization

- Advanced form of stemming
- reduces to dictionary base form
- example: "playing" → "play", "studies" → "study"
- requires POS(part of speech) tagging, slower but better

12.10 Handling spelling and special entities

- Spelling correction
- Handle emojis, mention, hashtags (social media)

12.11 Name Entity Recognition

- NER can be used for censoring sensity information
- for extracting NER and use it as features for the algorithms

12.12 Encoding

12.12.1 Bag of Words(BOW)

- Bag of words is one-hot encoding for words
- Dictionary based bag of words is often used for larger datasets
- It represent frequency of words in a vector format

12.12.2 Term frequency-Inverse document frequency (TF-IDF)

- Mesure How important a word is to a document relative to the whole corpus
- common words, less importance
- Rare document specific word, higher importance

- Vectorize a word BOW style and find TF-IDF for each words to form a vector which can be used alongside with cosine similarity to find similar documents
- Drawbacks - doesn't consider order of the words

$$TFIDF(t, d) = TF(t, d) \times IDF(t)$$

Term Frequency is defined as number of times term t appears in document d by total number terms in the document

$$TF(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

Inverse Document Frequency is a penalty for words that appear in many documents

$$IDF(t) = \log \left(\frac{N}{DF(t)} \right)$$

12.12.3 Embedding

- Word2vec, embedding models
- converts words/sentences to vectors
- the embedding encodes the meaning of the sentence / word
- Words similar to each other are near to each other

12.13 Image processing

12.13.1 Resizing & Cropping

Scaling and cropping to have fixed dimension either to ingest the images into a fixed pipeline

Example: ResNet expects a 224 X 224

12.13.2 Normalization/Scaling

Use Min Max normalization to normalize 0 - 255 to 0-1, prevents large gradients, helps faster training

12.13.3 Grayscale conversion

Converting RGB(3 channels) to a single channel (greyscale) reduces data complexity/size and noise especially when it doesn't matter to have all the channels like xray and so

12.13.4 Data Augmentation

Modifying the input data to simulate real world data like flipping/ rotating the images, cropping, zooming, adjust brightness, contrast to add noise, prevents overfitting and improves generalization

12.13.5 Noise Recution/Smoothing

Camera might introduce some random noise in pixels, using a gaussian blur smoothes the image and let's the model identify patterns instead of noise

12.13.6 Histogram Equalization

Histogram Equalization boosts contrast by spreading pixel intensities, it re-distributes the pixel intensities so that the histogram covers the whole range (0-255), used to color correct low light or medical diagnosis photos

12.13.7 Color Standardization/White balance

You can either shoot a white/grey photo in the subjects lighting to calculate white balance or assume average color of your image should be grey and correct each channels accordingly

12.13.8 Segmentation

Use pretrained segmentation models like Mask R-CNN, U-net, deeplap, to predict the subjects and create masks for it

13 Feature engineering & selection

13.1 Filter Methods

Filter methods uses statistics to determine wheather the feature is fit to be in the model

13.1.1 Correlation coefficient

- Measures how strongly linearly(Pearson) or montonically(Spearman, Kendall) related to the target
- Drop Features that are highly correlated with each other (redunadant features)
- Don't drop features that are less correlated with the target, since they could be correlated in a complex way

13.1.2 Chi square test

Chi square is like MSE for the model that assumes features are independent, if the error aka chi square is high then they must be dependent, calculate chi-square for target variable and feature, if they are low disgard it

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

where,

- O = Observed Frequency
- E = Expected Frequency

13.1.3 ANOVA F-test

ANOVA stands for Analysis of Variance. it's a statistical method used to compare three or more group means to see if at-least one is significantly different, if you got two groups use t-test

$$F = \frac{\text{Variance between groups}}{\text{How much data varies inside each group}}$$

and calculate p value based on the f-score, and degree of freedom of both numerator and denominator. with p value you can either reject or accept the null's hypothesis

Nulls Hypothesis H_0 - All groups are the same

Alternate Hypothesis H_a - one or more group is different from each other

13.1.4 Variance Threshold

Calculate variance of the data if variance is low, discard it (too little information)

13.2 Wrapper Methods

13.2.1 Forward selection

1. Start with empty feature
2. train on all features
3. evaluate the model (use metrics + cross validation)
4. select the best forming one
5. repeat 2 - 4 till you hit maximum number of features or performances changes are negligible or negative

13.2.2 Backward Elimination

Start with all features, train your model using full feature set, evaluate statistical significance (p-value) remove the one with highest p-value, refit the model, repeat until all features are statistically significant ($p < 0.05$)

1. Start with all features
2. Train the model with the feature set
3. Evaluate statistical significance (p-value)
4. Remove the one with highest p-value
5. Repeat 2 - 4 till all the features are statistically significant ($p < 0.05$)

Where,

- Null Hypothesis(H_0) - The features has no effect
- Alternate Hypothesis(H_a) - The feature has an effect

13.2.3 Recursive Feature Elimination(RFE)

1. Start with all features
2. Train the model with the feature set
3. Evaluate Feature importance
4. Remove the one with least feature importance
5. Repeat 2 - 4 till you hit the desired number of features

13.3 Embedded Methods

13.4 Dimensionality Reduction

13.4.1 Principal Component Analysis (PCA)

13.5 Domain Knowledge

Research about the target variable and features and through domain knowledge you could predict/ assume wheather the feature is related to the target variable or not

Example: for titanic dataset, gender and survived are related because we know females are evacuated first

14 Supervised Learning

15 Unsupervised Learning

16 Time Series

A sequence of data points collected at specific points at time

16.1 Properties of Time series

16.1.1 Order

The data points are arranged in specific sequence based on time

16.1.2 Time dependence

Observations are often dependent on previous data points(lags)

16.1.3 Stationarity

A time series is considered to be stationary if it's statistical properties doesn't change, constant mean, constant variance, constant covariance

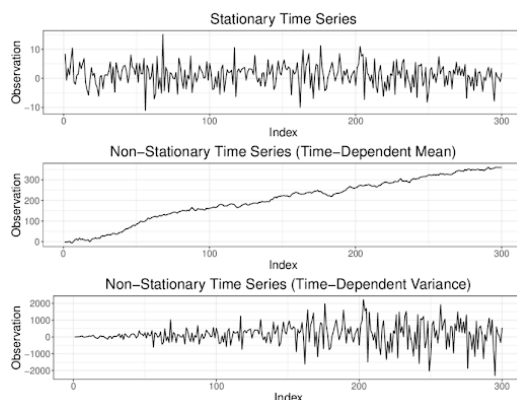


Figure 3: "Stationary Vs Non-stationary"

16.1.4 Trend

Long time increase or decrease in the data , could be linear exponential or other forms.

16.1.5 Seasonality

Regular, repeating patterns at fixed interval (daily, weekly, yearly)

Example: Ice cream sales peak every summer

16.1.6 Cyclic

Long term oscillations that doesn't have a fixed-period like Seasonality

Example: Business cycles - Booms and recessions

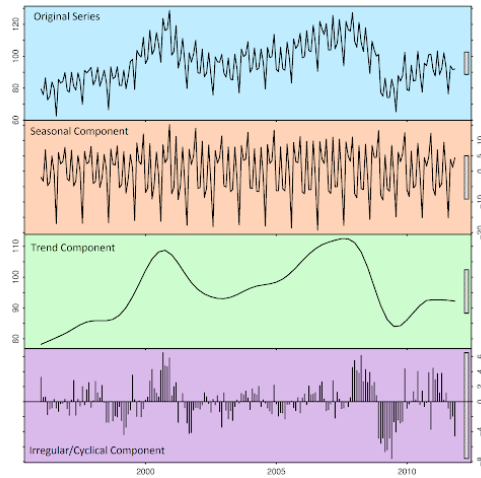


Figure 4: "Stationary Vs Non-stationary"

16.2 Suitability Test for TSA

How to find if traditional Time Series Analysis is suitable for the given data

16.2.1 Time dependent

The given data should have a temporal dependency(time), TSA is about finding patterns across time

16.2.2 Auto Correlation

Conventional time series models often rely on the assumption that current data point is dependent on previous data points(lags), ACF(auto correlation function) or PACF(Partial ACF) to check the auto correlation

16.2.3 Stationary

Many TSA models assumes stationarity, you can still model a non-stationary model by making it stationary, it should atleast be feasible to make it as a stationary time series

16.2.4 Data Quality

Data quality is important for any machine learning model, outliers, sufficient data, data completeness, etc...

16.2.5 White noise

Data shouldn't be a pure white noise (random), you can't predict a random data

White noise is a time series that satisfies following conditions

- Zero mean
- Constant Covariance
- Zero correlation with lags

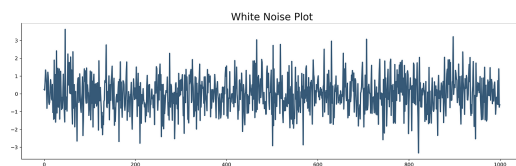


Figure 5: "Stationary Vs Non-stationary"

17 Tests for time series

Stationarity

17.1 Augmented Dicky Fuller Test (ADF)

Tests if a given time series is stationary or not, Ensures time series is stationary, critical for AR,MA,ARMA, and determines wheather the differencing is neccessary in ARIMA

The Hypothesis of the test is as follows,

- Null Hypothesis (H_0): Time series has a unit root - non-stationary
- Alternate Hypothesis (H_0): Time series does not have a unit root - stationary

Calculate p-value,

- if $p < 0.05$ reject null's hypothesis stationary
- if $p > 0.05$ failed reject null's hypothesis - non-stationary

ADF is often used to find d - number of differencing to make a TS stationary start with $d = 0$, increase until the null's hypothesis is rejected, usually you can stop after $d = 2$ as anything above 2 is super rare

17.2 Unit Root

To find unit root we assume the given time series unit root basically means the coefficient of a past term is 1 (unit), thus it's effects persists if it can be found using the test we can use differencing(to make it stationary) and pass it to ARMA model, basically ARIMA to negate the unit root effect, we can use dicky fuller test to determine how many differencing is required for the TS to be

Auto Corelation

17.3 Auto Correlation Function(ACF)

Auto correlation Function is a tool used to measure

17.4 PACF

Seasonality and Trend

17.5 Time series decomposition

17.6 Periodogram

Variance/ Volatility

17.7 ARCH test

18 Evaluation Metrics for Regression

18.1 Mean Absolute Error

$$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i|$$

- Robust to outliers, treats all errors equally doesn't square the errors like RMSE, MSE..etc
- It's used when your model can tolerate moderate outliers
- Interpretability - Has same unit as the thing you are predicting/easy to understand
- Gives out constant gradient (bad for gradient based loss function)

18.1.1 Gradient of MAE

$$\frac{d}{d\hat{y}}|y - \hat{y}| = \begin{cases} +1 & \text{if } \hat{y} < y \\ -1 & \text{if } \hat{y} > y \\ \text{undefined} & \text{if } \hat{y} = y \end{cases}$$

As you can see no matter how far the error is from true value it always gives a constant gradient as it treats every error as same stics

18.2 Mean Squared Error

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

- Penalizes large errors/outliers
- Gives out strong gradient signals

18.2.1 Gradient of MSE

$$\frac{dMSE}{d\hat{y}} = -\frac{2}{n}(y - \hat{y})$$

It points in the direction of the error, and it grows linearly with size of the error. Larger the gradient, when prediction are more wrong \rightarrow model adjusts faster

18.3 Root Mean Squared Error(RMSE)

$$RMSE = \sqrt{MSE}$$

- It combines interpretability of MAE and sensitive to errors of MSE
- It has smooth gradient curves just like MSE, and it's preferred for gradient descent

18.3.1 Gradient of RMSE

$$\frac{dRMSE}{d\hat{y}_i} = \frac{1}{n \times RMSE}(\hat{y}_i - y_i)$$

1. The gradient strength changes with RMSE, if your RMSE is very large the gradient becomes small, and if your RMSE is very small the gradient becomes large.
2. It makes RMSE a Non-constantly scaled loss
3. MSE is preferred over RMSE in training, but RMSE is preferred while reporting for interpretability

18.4 R-Square(R^2)

R^2 is the coefficient of determination. it tells how well your regression model explains the variation in the dependent variable(Y) using independent variables(X)

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

where,

- $SS_{res} = \sum (y_i - \hat{y}_i)^2 \rightarrow$ Residual sum of squares(error)/MSE
- $SS_{tot} = \sum (y_i - \bar{y}_i)^2 \rightarrow$ Total sum of squares (total variability)

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y}_i)^2}$$

Or, it can also be written more intuitively as

$$R^2 = 1 - \frac{MSE}{\sigma^2}$$

- Let us understand the formula (1-) operator just switches from maximizing to minimizing so you can ignore that.
- $\frac{MSE}{\sigma^2}$ Explains how well our model performs to a model that just predicts mean everytime, so if the ratio is 1, then our model is same as the dumb model, we have to reduce the ratio but the world likes "more the better" approach add (1-) operator we have to maximize the error and it's called as R^2
- R^2 ranges from $(-\infty, 1]$

18.5 Adjusted R^2

$$R_{adj}^2 = 1 - \left(\frac{(1 - R^2)(n - 1)}{n - k - 1} \right)$$

The above mentioned is textbook formula but we use our simplified representation for R^2

$$R^2 = 1 - \frac{MSE}{\sigma^2}$$

so, R_{adj}^2 would be

$$R_{adj}^2 = 1 - \frac{MSE_{adj}}{\sigma_{adj}^2}$$

MSE adjusted accounts for the number of freedoms used up to predict the data, which is K, represents number of parameters like number of predictors, number of bias

$$MSE_{adj} = \frac{\sum (y_i - \hat{y}_i)^2}{n - k}$$

Variance adjusted for number of freedoms used up which is 1 (mean), thus it'd be n-1 insted of n

$$\sigma_{adj}^2 = \frac{\sum (y_i - \mu)^2}{n - 1}$$

Substituting we get,

$$R_{adj}^2 = 1 - \left(\frac{MSE}{\sigma^2} \times \frac{n - 1}{n - k} \right)$$

where

- n - number of samples/ training samples
- k - number of parameters

18.6 Mean Absolute Percentage Error(MAPE)

MAPE is a metric used to measure accuracy of a predictive model. It expresses the prediction error as the percentage of actual values

$$MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100$$

MAPE is just like MAE but it gives out the error in percentage thus it's easier to intrepret

18.7 Huber Loss

Huber loss is a robust loss function/evaluation metric that has both strengths of MAE and MSE

$$L_{\delta} = \begin{cases} \frac{1}{2}(y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta \\ \delta \cdot (|y - \hat{y}| - 1/2\delta) & \text{otherwise} \end{cases}$$

where,

- δ is a hyperparameter that controls the behavior between MSE and MAE behavior

18.7.1 Gradient of Huberloss

$$\frac{\partial L}{\partial \hat{y}} = \begin{cases} -(y - \hat{y}) & \text{if } |y - \hat{y}| \leq \delta \\ -\delta \cdot \text{sign}(y - \hat{y}) & \text{otherwise} \end{cases}$$

example graph goes here

19 Evaluation Metrics for Classification

19.1 Basic Terminologies

19.1.1 True Positive

Correctly predicted positive class

19.1.2 False Positive

Falsely Predicted Positive class, actually negative

19.1.3 True Negative

Correctly predicted negative class

19.1.4 False Negative

Falsely predicted negative class, actually positive

19.2 Accuracy

Out of all predictions how many are correct, ratio between correct predictions and total predictions is accuracy

$$Accuracy = \frac{\text{Correct predictions}}{\text{Total predictions}}$$

can also be written as,

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

19.3 Precision

Out of all my positive class prediction how much did I get correctly, ratio between correct positive class prediction and total positive class predicted

$$\text{Precision} = \frac{\text{No of correct positive class predicted}}{\text{No of positive class predicted}}$$

can also be written as,

$$\text{Precision} = \frac{TP}{TP + FP}$$

Use when FP is costly, when you want every positive prediction to be trust worthy, measures reliability

19.4 Negative Predicted Value (NPV)

Out of all my negative predictions how much did I get correctly, ratio between

$$NPV = \frac{\text{No of correct negative class predicted}}{\text{No of negative class predicted}}$$

can also be written as,

$$NPV = \frac{TN}{TN + FN}$$

Use when FN is costly, when you want every negative class to be trust worthy

19.5 Recall(TPR)

Out of all positive cases how many did I predict correctly, also known as True Positive Rate(TPR) ratio of correctly predicted positive class and total positive cases

$$\text{Recall} = \frac{\text{No of correctly predicted positive class}}{\text{No of actual positive class}}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

Use when FN is costly/ wrongly classifying , did it recall everything, predict all positive cases

19.6 Specificity

Out of all negative cases how many did I predict correctly, Basically Recall for negative class

$$\text{Specificity} = \frac{\text{No of correctly predicted negative class}}{\text{No of actual negative class}}$$

$$\text{Specificity} = \frac{TN}{TN + FP}$$

Use when FP is costly, when predicting negative class is important

19.7 False Positive Rate(FPR)

Out of all negative cases how many did I fail to predict,

$$\text{FPR} = \frac{\text{No of wrongly predicted negative class}}{\text{No of actual Negative class}}$$

$$\text{FPR} = \frac{FP}{TN + FP}$$

$$\text{FPR} = 1 - \text{specificity}$$

19.8 F-score

It's a harmonic mean between precision and recall

$$F_{\beta} = \frac{(1 + \beta^2) \cdot \text{Precision} \cdot \text{Recall}}{(\beta^2 \cdot \text{Precision}) + \text{Recall}}$$

19.8.1 Derivation

$$\frac{1}{F_\beta} = \frac{w_p}{\text{Precision}} + \frac{w_r}{\text{Recall}}$$

We want harmonic ratio of precision and recall, β is how much more you care about recall then precision

$$\frac{w_r}{w_p} = \beta^2$$

We want the weights to add up to one

$$w_r + w_p = 1$$

By solving we get

$$w_r = \frac{\beta^2}{1 + \beta^2}$$
$$w_p = \frac{1}{1 + \beta^2}$$

20 AUC-ROC curve

20.1 Receiver Operating Characteristic (ROC)

The ROC curve plots the True Positive Rate(TPR) vs False Positive Rate(FPR) at various threshold settings, it let's you decide which one is best for you based on your requirement

20.2 Area Under the Curve (AUC)

The AUC tells how much of the curve is under the line, usually compared with other models

Higher AUC = Better model performance

Auc Score	Intrepretation
0.5	Random guessing
0.7 - 0.8	Acceptable
0.8 - 0.9	Excellent
>0.9	Outstanding

21 Regularization

Regularization adds a penalty for complexity to prevent overfitting, make model simpler

$$\text{Regularized Loss} = \text{Loss} + \text{Regularization parameter}$$

21.1 L1

L1 adds the sum of absolute value of coefficients to the loss function, it is used for feature selection since it encourages the optimizer to shrink unwanted features to zero

$$\text{L1} = \lambda \sum |w_i|$$

21.2 L2

L2 adds the sum of coefficient squares to the loss function this makes the gradient strength smooth, thus it won't reduce everything to zero, but towards zero, use it when you think all features contribute to ur model

$$\text{L2} = \lambda \sum w_i^2$$

21.3 Elastic net

It's a combination of L1 and L2 where you want to have both sparsity(l1) and stability(l2)

$$\text{Elastic Net} = \alpha \text{L1} + (1 - \alpha) \text{L2}$$

22 TODOO

1. pca
2. complete the Distributions, some intro to probability
3. Time series, ML, other stuff