



BITS Pilani

Session 8

Predictive Analytics - Correlation



Agenda

- Problems on Chi Square distribution
- Correlation
- Revision of the topics covered till session 7

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 354

Example:

consider the following data

Travel time	Stress			Total
	High	Moderate	Low	
< 20 min	9	5	18	32
20-50 min	17	8	28	53
≥ 50 min	18	6	7	31
Total	44	19	53	116

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 355

Based on this data, Can we conclude that stress levels depends on travel time

???

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 356

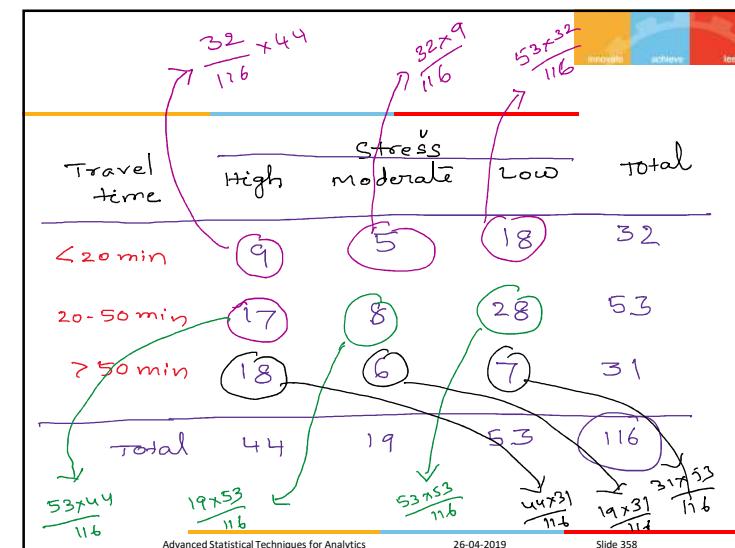
Chi-Square (χ^2) distribution

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

O: observed frequencies
E: expected frequencies

for $(r-1) + (c-1)$ degrees of freedom

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 357



$$\chi^2 = \sum \frac{(O-E)^2}{E} = \sum \frac{(9-12.14)^2}{12.14} \dots = 9.836$$

		Stress	Total	
Travel time	High	Moderate	Low	
< 20 min	9 / 12.14	5 / 5.24	18 / 14.62	32
20-50 min	17 / 20.10	8 / 8.68	28 / 24.22	53
≥ 50 min	18 / 11.75	6 / 5.08	7 / 4.17	31
Total	44	19	53	116

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 359

calculated $\chi^2 = 9.836$

Let $\alpha = 0.01$

d. o. f: $(r-1) \times (c-1)$
 $= (3-1) \times (3-1) = 4$

$$\chi^2_{0.01, 4} = 13.30$$

$\chi^2 = 9.836 < 13.30$

H_0 accepted

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 360

Example

A tobacco company claims that there is no relationship between smoking and lung ailments.

	Lung ailment	non-lung ailment	Total
Smokers	75	105	180
Non-smokers	25	95	120
Total	100	200	300

Based on this data, can we accept/reject the claim?

H_0

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 361

observed frequency

$$E = \frac{180}{300} \times 100 = 60$$

$$E = \frac{180}{300} \times 200 = 120$$

	Lung ailment	non-lung ailment	Total
Smokers	75	105	180
Non-smokers	25	95	120
Total	100	200	300

$$\frac{120}{300} \times 100 = 40$$

$$\frac{120}{300} \times 200 = 80$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 362

$$\chi^2 = \frac{(75-60)^2}{60} + \frac{(105-120)^2}{120} + \dots$$

From χ^2 -tables at 0.05 LOS
d.f. $(n-1) \times (c-1)$
 $(2-1) \times (2-1) = 1$

$$\chi^2 = 14.063 > 3.841$$

Reject H_0 →

	Lung ailment	non-lung ailment	Total
Smokers	75 (60)	105 (120)	180
Non-smokers	25 (40)	95 (80)	120
Total	100	200	300

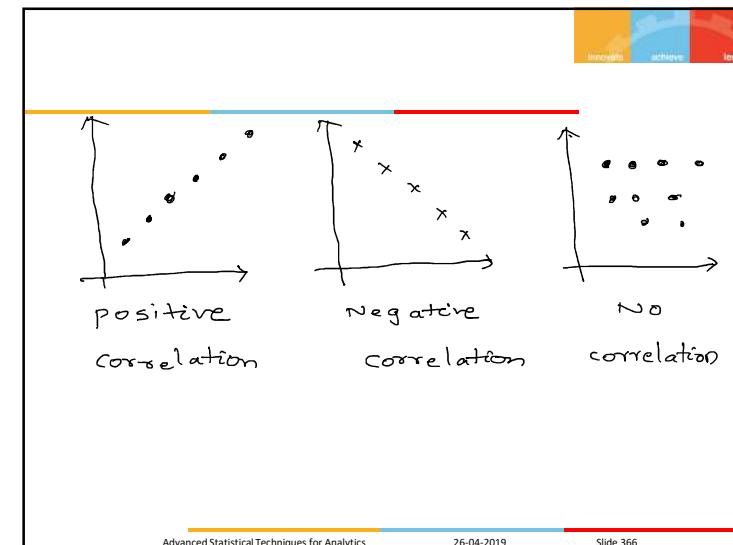
Advanced Statistical Techniques for Analytics 26-04-2019 Slide 363

Correlation

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 364

correlation

→ Sales of a company and Expenditure on advertisement
 → Price and Demand of a product
 → Inflation and Gold price
 → IQ and performance in Entrance.



Coefficient of Correlation

$r = 1 \Rightarrow$ perfect and positive relation
 $r = -1 \Rightarrow$ " " negative relation
 $r = 0 \Rightarrow$ no relation
 $0 < r < 1 \Rightarrow$ partial positive relation
 $-1 < r < 0 \Rightarrow$ " negative "

Coefficient of correlation:

$$r = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = \frac{\sum dx dy}{\sqrt{\sum dx^2 \sum dy^2}}$$

where $dx = x - \bar{x}$
 $dy = y - \bar{y}$
 $dx^2 = (x - \bar{x})^2$
 $dy^2 = (y - \bar{y})^2$

Example - 1

X	1	2	3	4	5	6	7	8	9
Y	10	11	12	14	13	15	16	12	18
\bar{x}	$\bar{x} = \frac{\sum x}{n} = \frac{45}{9} = 5$								
\bar{y}	$\bar{y} = \frac{\sum y}{n} = \frac{126}{9} = 14$								
r									

X	d_x	d_x^2	Y	d_y	d_y^2	$d_x d_y$
1	-4	16	10	-4	16	16
2	-3	9	11	-3	9	9
3	-2	4	12	-2	4	4
4	-1	1	14	0	0	0
5	0	0	13	-1	1	0
6	1	1	15	1	1	1
7	2	4	16	2	4	4
8	3	9	17	3	9	9
9	4	16	18	4	16	16

$r = \frac{\sum d_x d_y}{\sqrt{\sum d_x^2 \sum d_y^2}}$

$$= \frac{59}{\sqrt{60 \times 60}} = 0.9833$$

Coefficient of Determination

r is coeff. of correlation

r^2 is coeff of determination

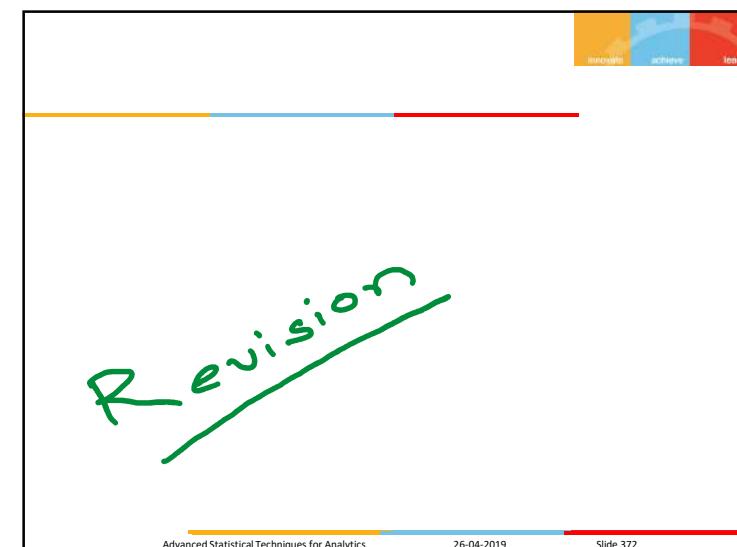
↓

Indicates the extent to which variation in one variable is explained by the variation in the other.

$r = 0.9 \Rightarrow r^2 = 0.81$

is 81% of the variation in y due to variation in x .

remaining 19% is due to some other factors



probability

1) Probability $P(A) = \frac{m}{n}$

2) $0 \leq P(A) \leq 1$

3) $P(A) + P(\bar{A}) = 1$

4) A and B are mutually exclusive events, then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 373

5) If A and B are independent events
 $\Leftrightarrow P(A \cap B) = P(A) \cdot P(B)$

6) Set theory
 happening of both $A \cap B$
 at least one ... $A \cup B$
 not happening --- \bar{A} etc

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 374

Conditional probability

7) $P(A|B) = \frac{P(A \cap B)}{P(B)}$

$P(B|A) = \frac{P(A \cap B)}{P(A)}$

8) $P(A \cap B) = P(A|B) P(B)$
 $= P(B|A) P(A)$

9) If A, B are indept. events,
 $P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \cdot P(B)}{P(B)} = P(A)$
 Similarly $P(B|A) = P(A)$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 375

Baye's theorem

\therefore Total prob. $P(B) = \sum_{i=1}^n P(B|A_i) P(A_i)$

$$P(A_i|B) = \frac{P(B|A_i) P(A_i)}{\sum_{i=1}^n P(B|A_i) P(A_i)}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 376

Random Variables

Discrete
 $X = 1, 2, 3, 4$
continuous
 $X \in (1, 3)$

prob. dist. fn
 $P(x)$
 $P(x)$

(i) $0 \leq P(x) \leq 1$
(ii) $\sum P(x) = 1$

prob. density function
 $f(x)$
 $f(x)$

(i) $0 \leq f(x) \leq 1$
(ii) $\int f(x) dx = 1$

Prob. distributions

- 1) Bimomial distribution.

$$P(x) = m c_x p^x q^{m-x}$$

$$m = 0, 1, 2, \dots, x$$
- 2) Poisson distribution

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$x = 0, 1, 2, \dots$$

3. Normal distribution

$$P(x_1 \leq x \leq x_2) = \frac{x_2 - \mu}{\sigma}$$

$$P(z_1 \leq z \leq z_2) = F(z_2) - F(z_1)$$

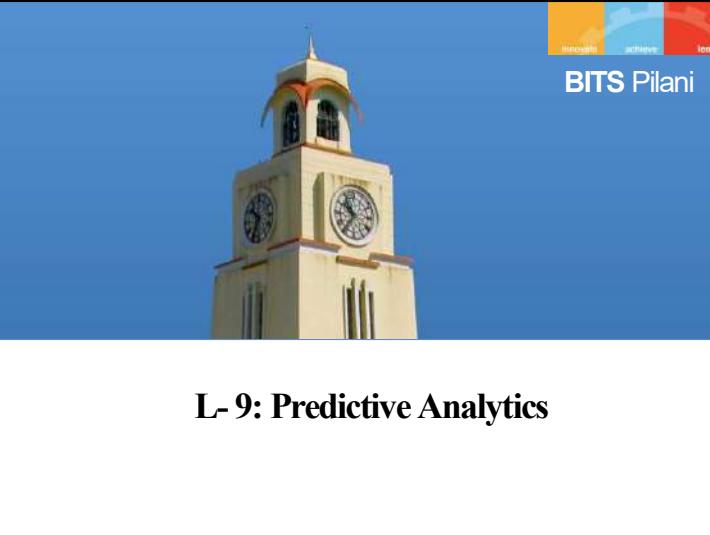
$F(z)$ from tables

$F(-z) = 1 - F(z)$

Testing of Hypothesis

```

graph TD
    A[Testing of Hypothesis] --> B[mean]
    A --> C[proportion]
    B --> D[one mean]
    B --> E[two means]
    C --> F[one]
    C --> G[two]
    C --> H[several]
    D --> I[Large]
    D --> J[small]
    I --> K[Z]
    J --> L[t]
    E --> M[Large]
    E --> N[small]
    M --> O[Z]
    N --> P[t]
    F --> Q[Z]
    G --> R[Z]
    H --> S[X2Dist]
  
```



L- 9: Predictive Analytics

Agenda

- Introduction to regression
- Method of least squares
- Simple linear regression
- Multiple linear regression

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 382

Regression

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 383

Regression :-

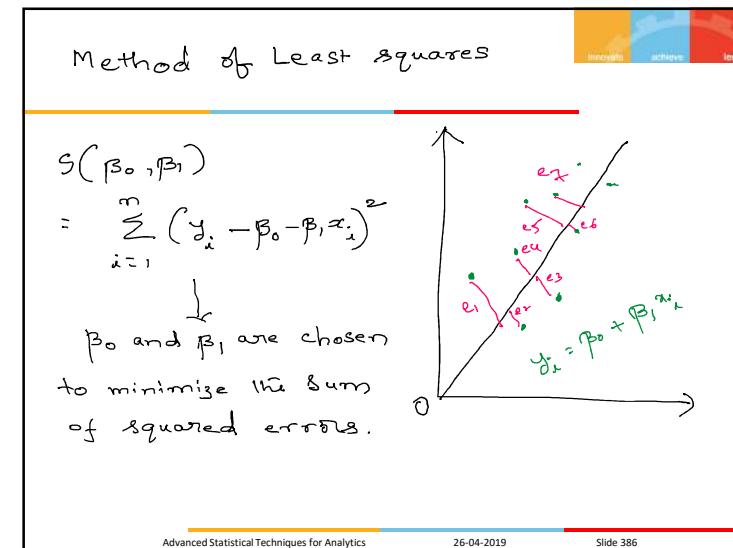
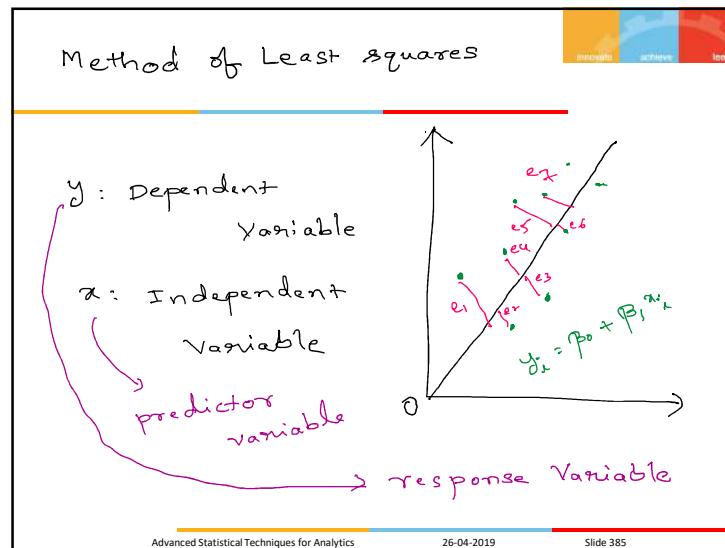
x	1	2	3	4	5
y	1	4	9	16	25

when $x = 7$: $y = ?$

x	1	2	3	4	5
y	1	6	2	5	4

when $x = 7$, $y = ?$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 384



Method of Least squares

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial S}{\partial \beta_0} = 0 \Rightarrow \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (-1)$$

$$\Rightarrow \sum_{i=1}^n y_i = n \beta_0 + \beta_1 \sum_{i=1}^n x_i$$

$$\frac{\partial S}{\partial \beta_1} = 0 \Rightarrow \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) (2)(-x_i)$$

$$\Rightarrow \sum_{i=1}^n x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2$$

on solving these, we get β_0 & β_1
which minimizes error.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 387

Linear regression

$$y = \beta_0 + \beta_1 x$$

$$\sum y = \beta_0 n + \beta_1 \sum x$$

$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

Normal equations.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 388

Matrix Approach:

Let $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{p-1} x_{p-1}$
 observations $y_i = 1, 2, \dots, n \rightarrow$ by a vector γ
 unknowns $\beta_0, \beta_1, \dots, \beta_{p-1} \rightarrow \beta$
 $x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np-1} \end{bmatrix}$
 $\therefore \hat{\gamma} = x \beta$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 389

Find β to minimize
 $s(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1 - \beta_2 x_2 - \dots)^2$
 $= \|y - x\beta\|^2 = \|y - \hat{y}\|^2$

Diffr S wrt to each β we get linear eqns

$$x^T x \hat{\beta} = x^T y \rightarrow \text{normal eqns}$$

If $x^T x$ is non-singular, the soln is

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 390

computationally, it is sometimes unwise even to form the normal equations because the multiplications involved in forming $x^T x$ can introduce undesirable round-off error.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 391

Linear regression (multiple regression)

Example:-

	size	no of rooms	no of floors	Age of home	price
1	2000	5	2	45	4000
1	1400	3	1	40	2000
1	1600	3	2	30	3000
1	800	2	1	35	2000

x_1 x_2 x_3 x_4 y

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 392

Example:

Consider the following data

x	1	2	4	0
y	0.5	1	2	0

Fit a linear regression line
Estimate y when x = 5.

$$y = \beta_0 + \beta_1 x$$

$$\sum y = n\beta_0 + \beta_1 \sum x$$

$$\sum xy = \beta_0 \sum x_1 + \beta_1 \sum x^2$$

$$3.5 = 4\beta_0 + \beta_1 (7)$$

$$10.5 = -7\beta_0 + \beta_1 (21)$$

on solving these

$$\beta_0 = 0$$

$$\beta_1 = 0.5$$

i.e. $y = 0 + (0.5)x$

When $x=5$, $y = (0.5)5$
 $= 0.25$



BITS Pilani

Session 10: Predictive Analytics(Continued)

Agenda

- Linear Regression(Revise)
- Model validation
- Ridge and lasso models
- Assumptions of Linear regression
- Logistic regression

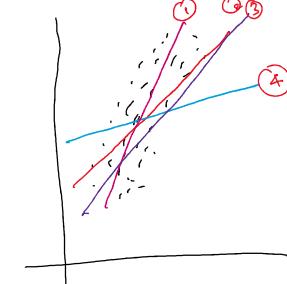
Classical Linear Regression (OLS)

- Explanatory and Response Variables are Numeric
- Relationship between the mean of the response variable and the level of the explanatory variable assumed to be approximately linear (straight line)
- Model:

$$Y = \beta_0 + \beta_1 x + \varepsilon \quad \varepsilon \sim N(0, \sigma)$$

Method of Least squares Errors

which one is the
best?



$$y = \beta_0 + \beta_1 x$$

$$\sum y = \beta_0 m + \beta_1 \sum x$$

$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

$$\downarrow \quad \beta_0 = ? , \beta_1 = ?$$

$$\text{then} \quad y = (\beta_0) + (\beta_1) x.$$

Multiple regression

Numeric Response variable (y)

p Numeric predictor variables

Model:

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon$$

innovate achieve lead
401/54

- Population Model for mean response:

$$E(Y | x_1, \dots, x_p) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- Least Squares Fitted (predicted) equation, minimizing SSE:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \hat{x}_1 + \dots + \hat{\beta}_p \hat{x}_p \quad SSE = \sum \left(Y - \hat{Y} \right)^2$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 401

innovate achieve lead
401/54

Accuracy of a model

By Using the following the strength of the linear model can be tested

- 1) Coefficient of determination (R^2)
- 2) Residual Standard error (RSE)

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 402

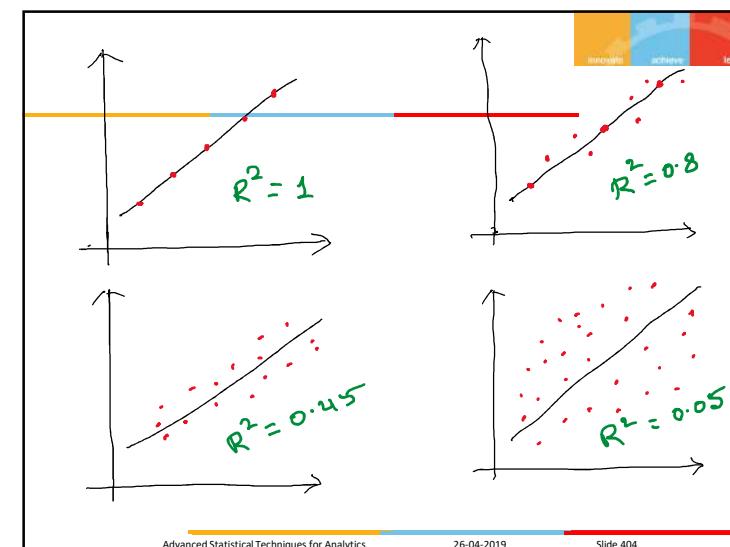
innovate achieve lead
401/54

$RSS \rightarrow$ Residual sum of squares
 $= \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$

$TSS \rightarrow$ $\sum_{i=1}^n (y_i - \bar{y})^2$ mean of respective variables

$$R^2 = 1 - \frac{RSS}{TSS}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 403



R – Squared vs Adjusted R - Squared

- In multiple regression, adjusted R – squared is better metric than R – squared asses the goodness of fit of the model
- R – squared always increases if additional variables are added into model , even if they are not related to the dependent variable

Regularization

- Over fitting can be solved with regularization
- Regularization can be done by putting constraints on the coefficients and variables.
- LASSO: Least Absolute Shrinkage and Selection Operator
Some coefficients can be dropped(i.e become zero)
- RIDGE: The coefficients will approach zero, but never dropped

Lasso & Ridge

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

• OLS estimation:

$$\min SSE = \sum \left(Y - \hat{Y} \right)^2$$

• LASSO estimation:

$$\min SSE = \sum_{i=1}^n \left(Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

• Ridge regression estimation:

$$\min SSE = \sum_{i=1}^n \left(Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^2$$

Assumptions in Regression Analysis

Assumptions

- The distribution of residuals is normal (at each value of the dependent variable).
- The variance of the residuals for every set of values for the independent variable is equal.
 - ✓ violation is called heteroscedasticity.
- The error term is additive
 - ✓ no interactions.
- At every value of the dependent variable the expected (mean) value of the residuals is zero
 - ✓ No non-linear relationships

409

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 409

- The expected correlation between residuals, for any two cases, is 0.
 - The independence assumption (lack of autocorrelation)
- ✓ All independent variables are uncorrelated with the error term.
- ✓ No independent variables are a perfect linear function of other independent variables (no perfect multicollinearity)
- ✓ The mean of the error term is zero.

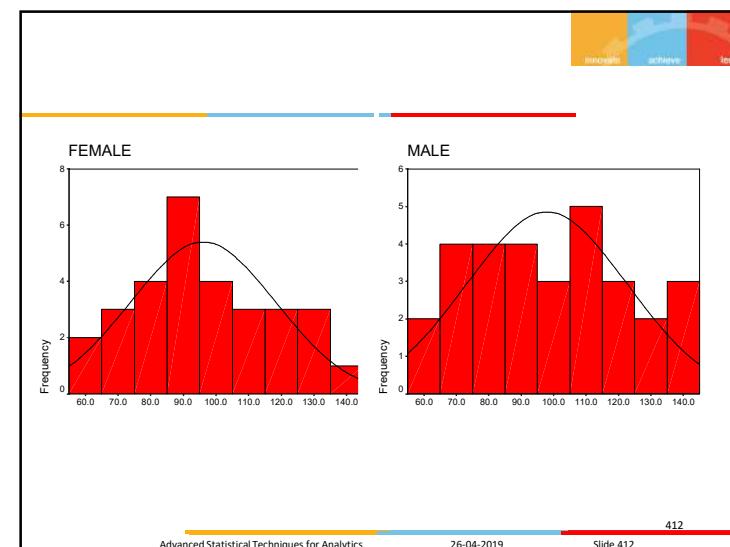
410

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 410

Assumption 1: The Distribution of Residuals is Normal at Every Value of the Dependent Variable

411

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 411



Non-Normality

Skew and Kurtosis

- Skew – much easier to deal with
- Kurtosis – less serious anyway

Transform data

- removes skew
- positive skew – log transform
- negative skew - square

Assumption 2: The variance of the residuals for every set of values for the independent variable is equal.

Heteroscedasticity

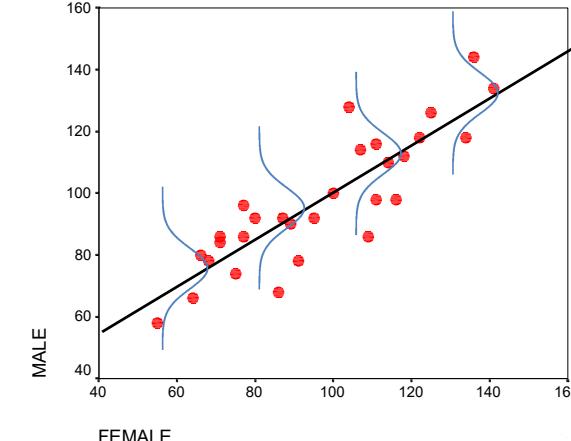
This assumption is about heteroscedasticity of the residuals

- Hetero=different
- Scedastic = scattered

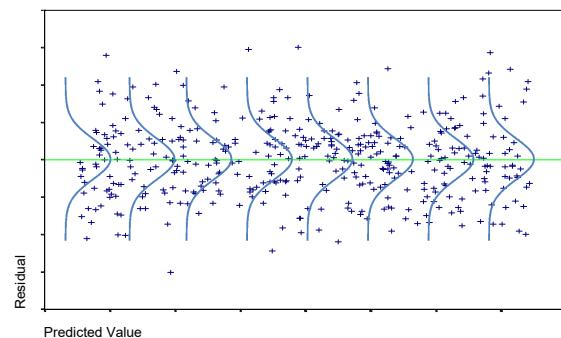
We don't want heteroscedasticity

- we want our data to be homoscedastic

Draw a scatterplot to investigate

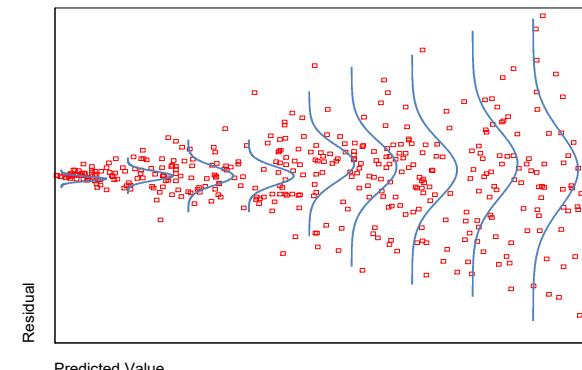


Good – no heteroscedasticity



417

Bad – heteroscedasticity



418

Assumption 3:
The Error Term is Additive

Assumption 4: At every value of the dependent variable the expected (mean) value of the residuals is zero

innovate achieve lead

Assumption 5: The expected correlation between residuals, for any two cases, is 0.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 421

•Result, with line of best fit

innovate achieve lead

Grade

Time

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 422

Now somewhat different

innovate achieve lead

Grade

Time

Question

- 3
- 2
- ◆ 1

423

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 423

innovate achieve lead

Assumption 6: All independent variables are uncorrelated with the error term.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 424

Assumption 7: No independent variables are a perfect linear function of other independent variables

Assumption 8: The mean of the error term is zero.

Multicollinearity

Correlation Matrix

	α_1	α_2	α_3	α_4
α_1	1	-0.80	0.98	0.061
α_2	-0.80	1	-0.184	0.103
α_3	0.98	-0.184	1	0.119
α_4	0.061	0.103	0.119	1

VIF(Variance Inflation Factor)

VIF(Variance Inflation Factor)

The better way to assess multi collinearity is to compute the VIF

$$\text{VIF} = \frac{1}{1 - R^2}$$

If VIF = 1 then Variables are not correlated

$1 < \text{VIF} < 5$ then the variables are moderately correlated

$\text{VIF} > 5$ then highly correlated and need to be eliminated from the model

innovate achieve lead

Logistic Regression

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 429

innovate achieve lead

Consider the following data

customer	Age	income (Lakhs)	gender	response
X	30	25	M	Yes 1
Y	45	40	F	No 0
Z	50	20	M	Yes 1
A	62	80	F	No 0
B	75	50	M	No 0
C	60	20	F	Yes 1

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 430

innovate achieve lead

we want to model this, then

$$y = \beta_0 + \beta_1 (\text{Age}) + \beta_2 (\text{income}) + \beta_3 (\text{gender})$$

Issues

1. Errors/residuals are not normally distributed.
2. No guarantee that the target/output/estimation is between 0 & 1.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 431

innovate achieve lead

Let us consider a function

$$f(z) = \frac{1}{1 + e^{-z}}$$

$$P(\text{response} = \text{yes} | \text{age}, \text{income}, \text{gender})$$

$$= \frac{1}{1 + e^{-(\beta_0 + \beta_1 \cdot \text{age} + \beta_2 \cdot \text{income} + \beta_3 \cdot \text{gender})}}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 432

Why use logistic regression?

- There are many important research topics for which the dependent variable is "limited."
- For example: voting, morbidity or mortality, and participation data is not continuous or distributed normally.
- Logistic regression is a type of regression analysis where the dependent variable is a dummy variable: coded 0 (did not vote) or 1(did vote)

Logistic Regression

Logistic regression is a supervised classification model. This allows us to make predictions from labelled data ,if the target variable is categorical.

Binary classification

Examples

1. A customer will default on a loan or not
2. A particular machine will break down in the next month or not
3. Predicting whether an incoming email is spam or not

Categorical Response Variables

Examples:

Whether or not a person smokes

Binary Response

Success of a medical treatment

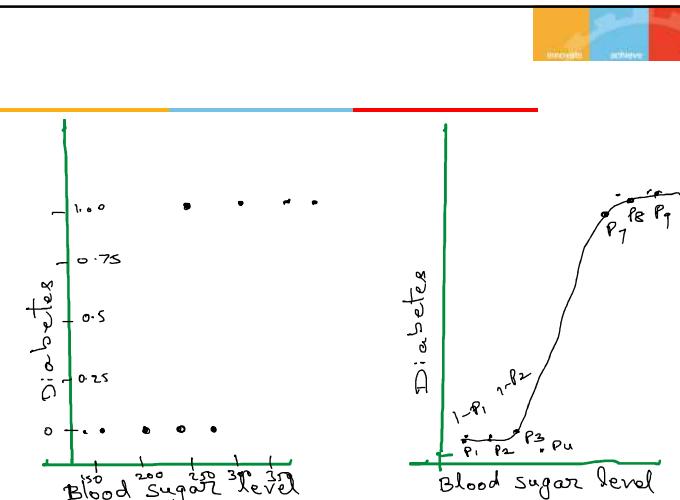
Opinion poll responses

Ordinal Response

$$Y = \begin{cases} \text{Non-smoker} \\ \text{Smoker} \end{cases}$$

$$Y = \begin{cases} \text{Survives} \\ \text{Dies} \end{cases}$$

$$Y = \begin{cases} \text{Agree} \\ \text{Neutral} \\ \text{Disagree} \end{cases}$$



innovate achieve lead

$$P(\text{diabetes}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Likelihood = $(1-p_1)(1-p_2)(1-p_3)(1-p_4)$
 $p_5(1-p_6)p_7p_8p_9p_{10}$

* $\left[(1-p_i)(1-p_j) \dots \text{for all non diabetics} \right]$.
* $\left[p_i \cdot p_i \dots \text{for all diabetes} \right]$.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 437

Y = Binary response X = Quantitative predictor

p = proportion of 1's (yes, success) at any X

Equivalent forms of the logistic regression model:

Logit form	Probability form
$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X$	$P = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$
	$\equiv \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}}$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 438

Binary Logistic Regression Model

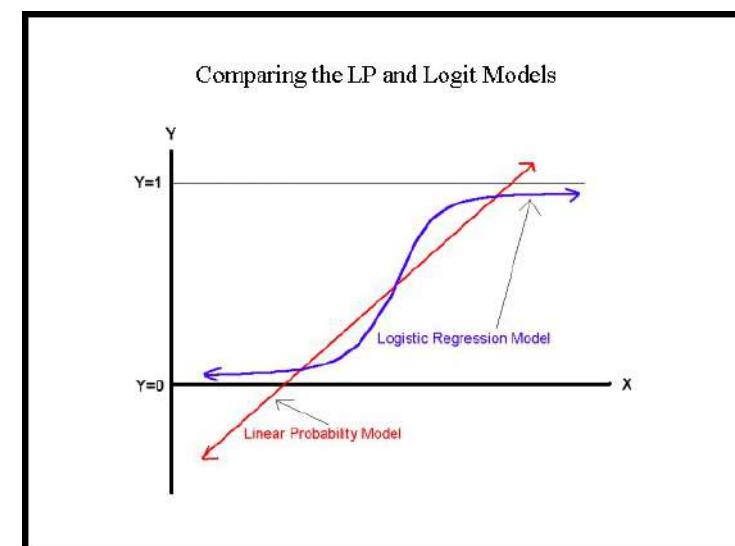
$Y = \text{Binary}$ $X_1, X_2, \dots, X_k = \text{Multiple}$

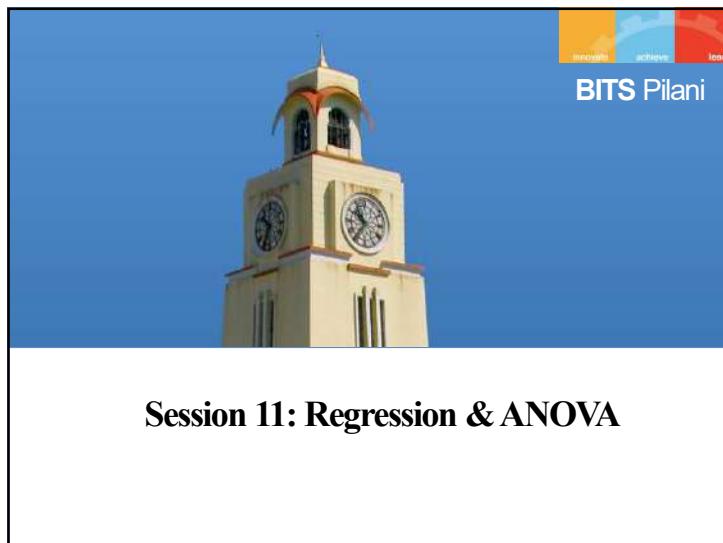
$\pi = \text{proportion of 1's at any } x_1, x_2, \dots, x_k$

Equivalent forms of the logistic regression model:

Logit form $\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$

Probability form $P = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k}}$





Session 11: Regression & ANOVA

Agenda

- Linear Regression(Revise)
- Model validation
- Ridge and lasso models
- Assumptions of Linear regression
- Logistic regression

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 442

Example:

consider the following data

y	370	386	443	499	528	616
x	22	25	28	31	33	38

Fit a regression line $y = \beta_0 + \beta_1 x$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 443

$$y = \beta_0 + \beta_1 x$$

$$\sum y = \beta_0 n + \beta_1 \sum x$$

$$\sum xy = \beta_0 \sum x + \beta_1 \sum x^2$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 444

x	y	x^2	xy
22	370	484	8140
25	386	625	9650
28	443	784	12404
31	499	961	15469
33	528	1089	17424
38	616	1444	23408
Σx	Σy	Σx^2	Σxy
177	2842	5387	86495

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 445

$$y = \beta_0 + \beta_1 x$$

$$\Sigma y = \beta_0 n + \beta_1 \Sigma x$$

$$\Sigma xy = \beta_0 \Sigma x + \beta_1 \Sigma x^2$$

i.e. $2842 = 6(\beta_0) + \beta_1 (177)$

$$86495 = 177 (\beta_0) + \beta_1 (5387)$$

on solving $\beta_0 = 0.195$

$$\beta_1 = 16.05$$

$$\therefore y = (0.195) + (16.05)x$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 446

Another form

$$(y - \bar{y}) = \beta_{xy} (x - \bar{x})$$

mean of y mean of x

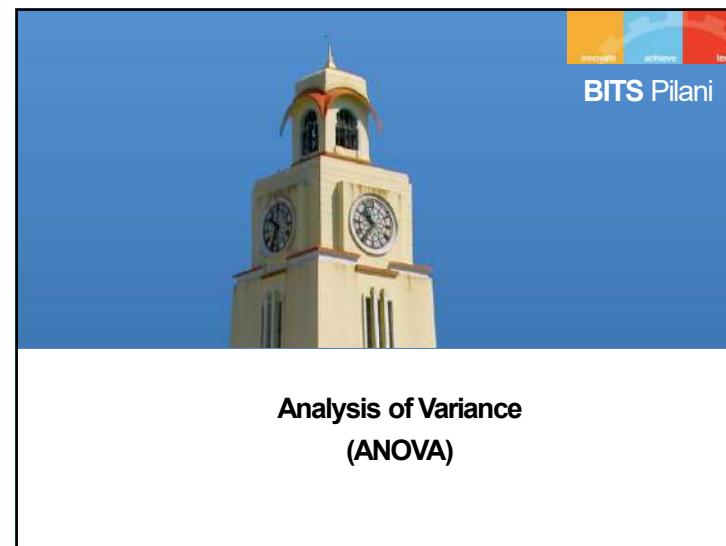
Regression coefficient of y on x

$$\beta_{xy} = \frac{\sigma_{xy}}{\sigma_x^2}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 447



• For example: A salesman in car sales wants to find the difference between two types of cars in terms of mileage:

- Mid-Size Vehicles



- Sports Utility Vehicles



Advanced Statistical Techniques for Analytics 26-04-2019 Slide 449

• For example: A salesman in car sales wants to find the difference between three types of cars in terms of mileage:

- Mid-Size Vehicles



- Sports Utility Vehicles



- PICKUP TRUCK



Advanced Statistical Techniques for Analytics 26-04-2019 Slide 450

Case

Testing the impact of nutrition and exercise on 60 candidates between age 18 and 50. They are grouped with different strategies. Now we need to find the most effective strategy

Group 1 eats only junk food

Group 2 eats only healthy food

Group 3 eats junk food & does cardio exercise every other day

Group 4 eats healthy food & does cardio

Group 5 eats junk food & does both cardio & strength training every other day

Group 6 eats healthy food.....

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 451

ANOVA

- ✓ Effectiveness of different promotional activities
- ✓ Quality of a product produced by different manufacturers in terms of an attribute
- ✓ Yield of crop due to varieties of seeds , fertilisers and quality of soil

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 452

ANOVA

- ✓ In real life things do not typically result in two groups being compared
- ✓ Used with 3 or more groups to test for MEAN DIFFS.
- ✓ We have at least 3 means to test, e.g., $H_0: \mu_1 = \mu_2 = \mu_3$.
- ✓ Could take them 2 at a time, but really want to test all 3 (or more) at once.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 453

Between/Within Groups

Variance can be separated into two major components

- **Within groups** – variability or differences in particular groups (individual differences)
- **Between groups** - differences depending what group one is in or what treatment is received

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 454

Null hypothesis:
no difference in means

Alternative hypothesis:
difference in means

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 455

ANOVA

Effectiveness of different promotional activities

Quality of a product produced by different manufacturers in terms of an attribute

Yield of crop due to varieties of seeds , fertilisers and quality of soil

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 456

ANOVA-analysis of variance

* Significance of difference between two sample means

$$H_0 = \mu_1 = \mu_2 = \dots = \mu_K$$

$$H_1 = \mu_1 \neq \mu_2 \neq \dots \neq \mu_K$$

Alternative Hypothesis

Assumptions

Each population is normally distributed with mean μ_i With equal variances

$$\sigma^2_i$$

Each sample is drawn randomly and independent of other samples

Observations	population		
	A	B	C
1	26	18	23
2	25	16	19
3	28	17	26
4	12	8	30

ANOVA summary

Source of Variation	Sum of squares	d.o.f	Mean squares	F-value
Between Samples	SSTR	n-1	MSTR = SSTR / (n-1)	
within Samples (Error)	SSE	n-n	MSE = SSE / (n-n)	$F = \frac{MSTR}{MSE}$
Total	SST	n-1		

Short cut method

$$T = \sum x_1 + \sum x_2 + \dots + \sum x_n$$

Corr. Fact : CF = $\frac{T^2}{n}$; $n = n_1 + n_2 + \dots + n_r$

$$ST = \left[\sum (x_1^2) + \sum (x_2^2) + \dots + \sum (x_n^2) \right] - CF$$

$$SSTR = \frac{\left(\sum x_j \right)^2}{n_j} - CF$$

$$SSE = SST - SSTR$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 461

Example

To test the significance of variation in the retail prices of a commodity in three metro cities, Mumbai, Kolkata and Delhi, four shops are chosen at random and the prices are given below

Mumbai : 16 8 12 14

Kolkata : 14 10 10 6

Delhi : 4 10 8 8

Prices in 3 cities are significantly different?

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 462

Example

A study was conducted to investigate the perception of corporate ethical values among individuals specialising in marketing. Using 0.05 level of significance and the data given below, test for significant differences in perception among three groups. (higher scores indicate higher ethical values)

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 463

	Marketing manager	Marketing Research	Advertising
1	6	5	6
2	5	5	7
3	4	4	6
4	5	4	5
5	6	5	6
6	4	4	6

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 464

$n=3, m=18$

$$T = \sum x_1 + \sum x_2 + \sum x_3$$

$$= 30 + 27 + 36 = 93$$

$$CF = \frac{T^2}{m} = \frac{(93)^2}{18} = 480.50$$

$$SST = (\sum x_1^2 + \sum x_2^2 + \sum x_3^2) - CF$$

$$= 154 + 123 + 218 = 495.50$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 465

$$SSTR = \left(\frac{\sum x_1^2}{m_1} + \frac{\sum x_2^2}{m_2} + \frac{\sum x_3^2}{m_3} \right) - CF$$

$$= \frac{30^2}{6} + \frac{27^2}{6} + \frac{36^2}{6} - 480.50$$

$$\approx 7$$

$$SSE = SST - SSTR$$

$$= 495.50 - 7 = 488.50$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 466

$$MSTR = \frac{SSTR}{df_1} = \frac{7}{2} : 3.5$$

$$MSE = \frac{SSE}{df_2} = \frac{7.5}{df_2} : 0.5$$

$$F = \frac{MSTR}{MSE} = \frac{3.5}{0.5} = 7$$

F-distribution tables

calculated value: 7
F-table value: 3.68

$3.68 < 7 \Rightarrow \text{Rejected}$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 467

Two way ANOVA

Sources of variation	Sum of square	Dof	mean square	test statistic
Between columns	SSTR	C-1	MSTR = $SSTR / (C-1)$	<i>* Treatment</i>
Between rows	SSR	n-1	MSR = $SSR / (n-1)$	<i>* Interaction</i>
Residual error	SSE	(C-1) * (n-1)	MSE = $SSE / ((C-1)(n-1))$	<i>* Error</i>
Total	SST	n-1		

** test statistic*

$MSTR = SSTR / (C-1)$

$MSR = SSR / (n-1)$

$MSE = SSE / ((C-1)(n-1))$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 468

Example

Analyse per year

Month	Salesman			
	A	B	C	D
May	50	40	48	39
June	46	48	50	45
July	39	44	40	39

→ Is there any significant difference in the sales made, etc.,
→ Is there a significant diff in the sales made during these months.

Salesman

	<i>x₁</i>	<i>x₂</i>	<i>x₃</i>	<i>x₄</i>
May	50	40	48	39
June	46	48	50	45
July	39	44	40	39

Ex: 15, 12, 18, 3

$$T = 15 + 12 + 18 + 3 = 48$$

$$CF = \frac{T^2}{n} = \frac{48^2}{12} = 192$$

SSTR = Sum of squares (columns)

$$= \left(\frac{15^2}{3} + \frac{12^2}{3} + \frac{18^2}{3} + \frac{3^2}{3} \right) - 192$$

$$= 42$$

SSR = Sum of squares between months (rows)

$$= \left(\frac{17^2}{4} + \frac{29^2}{4} + \frac{22^2}{4} \right) - 192$$

$$= 91.5$$

$$SST = \left(\sum x_1^2 + \sum x_2^2 + \sum x_3^2 + \sum x_4^2 \right) - CF$$

$$= (137 + 80 + 164 + 27) - 192$$

$$= 216$$

innovate achieve lead

$$\begin{aligned}
 SSE &= SST - (SSTR + SSR) \\
 &= 216 - (42.0 + 91.5) \\
 &= 82.5 \\
 df_c &= 3, \quad df_{n-1} = n-1 = 3-1 = 2 \\
 df_{(c-1)(n-1)} &= 3 \times 2 = 6
 \end{aligned}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 473

innovate achieve lead

$$\begin{aligned}
 MSTR &= \frac{SSTR}{c-1} = \frac{42.0}{3} = 14 \\
 MSR &= \frac{SSR}{n-1} = \frac{91.5}{2} = 45.75 \\
 MSE &= \frac{SSE}{(c-1)(n-1)} = \frac{82.5}{6} = 13.75
 \end{aligned}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 474

innovate achieve lead

	Sum of squares	D.o.f	mean squares	Variance ratio	
* Between Salesmen	SSTR 42.0	c-1 3	MSTR $\frac{42.0}{3} = 14$	$F_{\text{Treatment}}$ $\frac{14}{13.75} = 1.018$	$F > F_{\text{MSR}}$
* Between months	SSR 91.5	n-1 2	MSR 45.75	F_{block} $\frac{45.75}{13.75} = 3.327$	$F > F_{\text{MSR}}$
* residual error	SSE 82.5	(c-1)(n-1) 6	MSE 13.75		
Total	216	11			

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 475

innovate achieve lead

- $F_{\text{Treatment}} = 1.018 <$
 $\downarrow \text{columns}$
 $df_1 = 3, df_2 = 6$
 $\alpha = 0.05$
 accept
- $F_{\text{block}} = 3.327 <$
 $\downarrow \text{rows}$
 $2, 6$
 accept
 difference in the sales by salesmen
 difference in sales made during months.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 476



L- 12:Applied Multivariate Analytics

Agenda

- Bi Variate distribution
- multivariate normal distribution
- Eigen values and eigen vectors
- Principal Component Analysis

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 478

Joint Discrete Distribution Function _Discrete

$P(x_i, y_j) \rightarrow$ joint probability distribution function

$P_x(x_i, y_j) \rightarrow$ marginal prob. dist. fun of x

$P_y(x_i, y_j) \rightarrow$ " " " of y .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 479

Joint Discrete Distribution Function _Continuous

$f(x, y) \rightarrow$ joint p.d.f

$f_x(x, y) \rightarrow$ Marginal pdf of x

$f_y(x, y) \rightarrow$ Marginal p.d.f of y

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 480

Example:-

$$f(x,y) = \begin{cases} Kxy; & 0 < x < 4, 1 < y < 5 \\ 0; & \text{otherwise} \end{cases}$$

- 1) Find 'K' value?
- 2) $P(1 < x < 2, 2 < y < 3)$
- 3) $P(x \geq 3, y \leq 2)$
- 4) $f_x(x, y)$
- 5) $f_y(x, y)$

$$1) \int_{x=0}^4 \int_{y=1}^5 f(x,y) dy dx = 1$$

$$\begin{aligned} &\Rightarrow K \int_{x=0}^4 \int_{y=1}^5 xy dy dx = 1 \\ &\Rightarrow K \left[\frac{x^2}{2} \right]_{x=0}^4 \left[\frac{y^2}{2} \right]_{y=1}^5 = 1 \\ &\Rightarrow K \left[\frac{16}{2} - 0 \right] \left[\frac{25}{2} - \frac{1}{2} \right] = 1 \\ &\Rightarrow K = \frac{1}{96} \end{aligned}$$

$$2) P(x \geq 3, y \leq 2)$$

$$\begin{aligned} &= \int_{x=3}^4 \int_{y=1}^2 Kxy dy dx \\ &= \frac{1}{96} \left[\frac{x^2}{2} \right]_{x=3}^4 \left[\frac{y^2}{2} \right]_{y=1}^2 \\ &= \frac{1}{96 \times 2 \times 2} \left[4^2 - 3^2 \right] \left[2^2 - 1^2 \right] \\ &= \frac{7}{128} \end{aligned}$$

$$3) P(1 < x < 2, 2 < y < 3)$$

$$\begin{aligned} &= \int_{x=1}^2 \int_{y=2}^3 Kxy dy dx \\ &= \frac{1}{96} \left[\frac{x^2}{2} \right]_{x=1}^2 \left[\frac{y^2}{2} \right]_{y=2}^3 \\ &= \frac{5}{128} \end{aligned}$$

4) Marginal of X

$$f_x(x, y) = \int_{y=1}^5 f(x, y) dy$$

$$= \int_{y=1}^5 kxy dy = \frac{ky}{96} x \left[\frac{y^2}{2} \right]_1^5$$

$$= \frac{x}{192} (25 - 1) = \frac{x}{8}$$

$$\therefore f_x(x, y) = \frac{x}{8}, 0 < x < 4$$

5) Marginal of Y .

$$f_y(x, y) = \int_{x=0}^4 kxy dx$$

$$= \frac{ky}{96} \left[\frac{x^2}{2} \right]_0^4$$

$$= \frac{4y}{12}$$

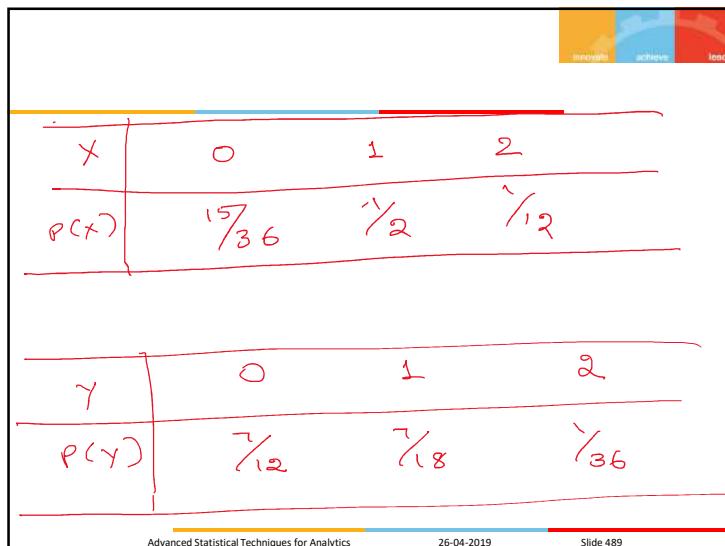
$$f_y(x, y) = \frac{y}{12}, 1 < y < 5$$

Example:

X	0	1	2
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{12}$
1	$\frac{2}{9}$	$\frac{1}{6}$	0
2	$\frac{1}{36}$	0	0

Example:

X	0	1	2	Marg. of Y
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{12}$	$\frac{1}{12}$
1	$\frac{2}{9}$	$\frac{1}{6}$	0	$\frac{1}{18}$
2	$\frac{1}{36}$	0	0	$\frac{1}{36}$
x	$\frac{15}{36}$	$\frac{1}{2}$	$\frac{1}{12}$	



Conditional Probabilities

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(x|y) = \frac{P(x,y)}{P_y(x,y)} \rightarrow \text{joint} \quad \text{Marginal of } Y$$

$$P(y|x) = \frac{P(x,y)}{P_x(x,y)} \rightarrow \text{joint} \quad \text{Marginal of } X$$

$$f(x|y) = \frac{f(x,y)}{f_y(x,y)} \rightarrow \text{joint} \quad \text{marginal of } y$$

$$f(y|x) = \frac{f(x,y)}{f_x(x,y)} \rightarrow \text{joint} \quad \text{marginal of } x$$

Independent r.v.s.

$$P(x,y) = P_x(x,y) \cdot P_y(x,y)$$

$$f(x,y) = f_x(x,y) \cdot f_y(x,y)$$

Bi-Variate Normal dist

$$f(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{x_1 - \mu_1}{2(1-\rho^2)}\right]$$

$$\chi^2 = \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} - 2\rho \frac{(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2}$$

$$\rho = \text{corr}(x_1, x_2) = \frac{\text{cov}(x_1, x_2)}{\sigma_1\sigma_2}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 493

Multivariate Normal distribution

$$\phi(x) = \left(\frac{1}{2\pi}\right)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}} \exp\left\{-\frac{1}{2}x^T \Sigma^{-1} x\right\}$$

x_1, x_2, \dots

Σ Det of Variance Covariance matrix

Σ^{-1} Inverse of Variance Covariance matrix

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 494

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Covariance between x_1 and x_2 .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 495

Preliminaries

- Standard Deviation is a measure of the spread of the data
- Variance – measure of the deviation from the mean for points in one dimension e.g. heights
- Covariance as a measure of how much each of the dimensions vary from the mean with respect to each other.
- Covariance is measured between 2 dimensions to see if there is a relationship between the 2 dimensions e.g. number of hours studied & marks obtained
- The covariance between one dimension and itself is the variance

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{cov}(x, x) = \text{variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 496

Covariance Matrix

- Representing Covariance between dimensions as a matrix
- e.g.

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix}$$

- Diagonal is the variances of x, y and z
- $\text{cov}(x,y) = \text{cov}(y,x)$ hence matrix is symmetrical about the diagonal
- N-dimensional data will result in nxn covariance matrix

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 497

- A positive value of covariance indicates both dimensions increase or decrease together
- A negative value indicates while one increases the other decreases, or vice-versa
- If covariance is zero: the two dimensions are independent of each other .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 498

Transformation matrices

Consider: $\begin{pmatrix} A \\ \text{---} \end{pmatrix} \times \begin{pmatrix} x \\ \text{---} \end{pmatrix} = \lambda \times \begin{pmatrix} x \\ \text{---} \end{pmatrix}$

$$\begin{pmatrix} 2 & 3 \\ 2 & 1 \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 12 \\ 8 \end{pmatrix} = 4 \times \begin{pmatrix} 3 \\ 2 \end{pmatrix}$$

Square transformation matrix transforms (3,2) from its original location. Now if we were to take a multiple of (3,2)

$$\begin{pmatrix} 2 \\ \text{---} \end{pmatrix} \times \begin{pmatrix} 3 \\ 2 \end{pmatrix} = \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

$$\begin{pmatrix} 2 \\ \text{---} \end{pmatrix} \times \begin{pmatrix} 6 \\ 4 \end{pmatrix} = \begin{pmatrix} 24 \\ 16 \end{pmatrix} = 4 \times \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 499

eigenvalue problem

- The eigenvalue problem is any problem having the following form:

$$A \cdot X = \lambda \cdot X$$

- A: n x n matrix
- X: n x 1 non-zero vector
- λ : scalar
- Any value of λ for which this equation has a solution is called the eigenvalue of A and vector v which corresponds to this value is called the eigenvector of A.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 500

eigenvalue problem

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix} = 4 \times \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$A \cdot X = \lambda X$

Therefore, $(3,2)$ is an eigenvector of the square matrix A and 4 is an eigenvalue of A

Given matrix A , how can we calculate the eigenvector and eigenvalues for A ?

eigenvalue *eigen vector*

$$AX = \lambda X$$

$$AX - \lambda I X = 0$$

$$[A - \lambda I] X = 0$$

Identity matrix
 $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

Homogeneous systems of equations

Non trivial soln (non-zero soln)

$$\begin{cases} 2x + 3y = 0 \\ 3x + 2y = 0 \end{cases}$$

iff

$$|A - \lambda I| = 0$$

Non trivial solns

Ex: $A = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$

$$|A - \lambda I| = 0 \Rightarrow \begin{vmatrix} 0-\lambda & 1 \\ -2 & -3-\lambda \end{vmatrix} = 0$$

$$\Rightarrow -\lambda(-3-\lambda) + 2 = 0$$

$$\Rightarrow 2 + 3\lambda + 2 = 0 \quad \therefore \lambda = -1, -2$$

Now we need to find X corresponding to $\lambda = -1$ and $\lambda = -2$ such that

$$AX = \lambda X$$

eigen values

when $\lambda = -1$: $[A - \lambda I] X = 0$

$$\therefore \begin{bmatrix} 0+1 & 1 \\ -2 & -2 \end{bmatrix} X = 0 \quad X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

non-trivial

$$\therefore x_1 + x_2 = 0$$

$$\text{and } -2x_1 - 2x_2 = 0 \quad \Rightarrow x_1 + x_2 = 0$$

Let $x_1 = K, x_2 = -K$

$$\therefore X_1 = \begin{bmatrix} K \\ -K \end{bmatrix} \quad \text{eigen vector}$$

Similarly we find eigen vectors corresponding to $\lambda = -2$

Data Presentation

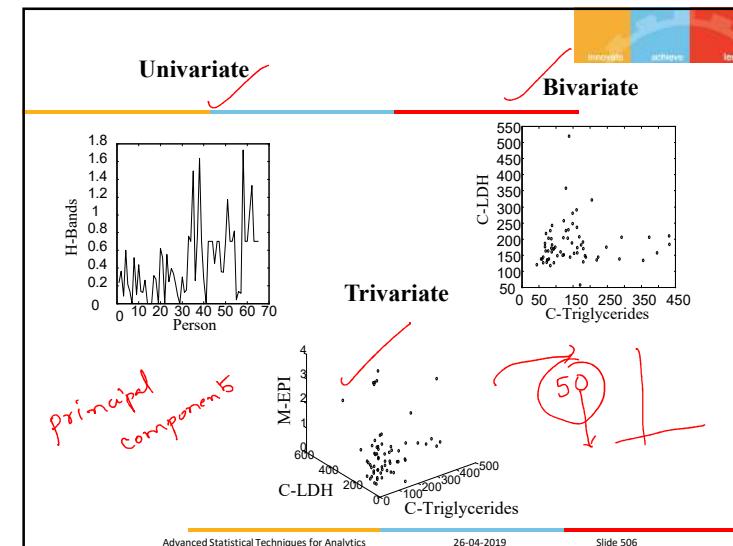
Blood and urine measurements (wet chemistry) from 65 people (33 alcoholics, 32 non-alcoholics).

Matrix Format

	HtMBC	HtRBC	HtHb	HtHd	HtMCV	HtMCH	HtMCHC
A1	8.0000	4.2000	14.7000	41.0000	85.0000	29.0000	34.0000
A2	7.3000	5.0200	14.7000	43.0000	98.0000	29.0000	34.0000
A3	4.3000	4.4800	14.1000	41.0000	91.0000	32.0000	35.0000
A4	7.5000	4.4700	14.9000	45.0000	101.0000	33.0000	33.0000
A5	7.3000	5.5200	15.4000	46.0000	84.0000	28.0000	33.0000
A6	6.9000	4.8800	16.0000	47.0000	97.0000	33.0000	34.0000
A7	7.8000	4.6800	14.7000	43.0000	92.0000	31.0000	34.0000
A8	8.6000	4.8200	15.8000	42.0000	88.0000	33.0000	37.0000
A9	5.1000	4.7100	14.0000	43.0000	92.0000	30.0000	32.0000

visualization

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 505



Applications

- Face Recognition
- Image Compression
- Gene Expression Analysis
- Data Reduction
- Data Classification
- Trend Analysis
- Factor Analysis
- Noise Reduction

100+

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 507

Principal Component Analysis

In real world data analysis tasks we analyze complex data i.e. multi dimensional data. We plot the data and find various patterns in it or use it to train some machine learning models. One way to think about dimensions is that suppose you have a data point x , if we consider this data point as a physical object then dimensions are merely a basis of view, like where is the data located when it is observed from horizontal axis or vertical axis.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 508

The header features three horizontal bars in orange, blue, and red, each containing a single word: "innovate", "achieve", and "lead".

As the dimensions of data increases, the difficulty to visualize it and perform computations on it also increases. So, how to reduce the dimensions of a data-

- * Remove the redundant dimensions
- * Only keep the most important dimensions

 DataAnalytics achieve

Now lets think about the requirement of data analysis. Since we try to find the patterns among the data sets so we want the data to be spread out across each dimension. Also, we want the dimensions to be independent. Such that if data has high covariance when represented in some n number of dimensions then we replace those dimensions with *linear combination* of those n dimensions. Now that data will only be dependent on linear combination of those related n dimensions. (*related = have high covariance*)



innovate achieve lead

- It is a linear transformation that chooses a new coordinate system for the data set such that
 - greatest variance by any projection of the data set comes to lie on the first axis (then called the first principal component),
 - the second greatest variance on the second axis, and so on.
- PCA can be used for reducing dimensionality by eliminating the later principal components.



innovate achieve le

what does Principal Component Analysis (PCA) do?

PCA finds a new set of dimensions (or a set of basis of views) such that all the dimensions are orthogonal (and hence linearly independent) and ranked according to the variance of data along them. It means more important principle axis occurs first. (more important = more variance/more spread out data)

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 512

How does PCA work

- Calculate the covariance matrix X of data points.
- Calculate eigen vectors and corresponding eigen values.
- Sort the eigen vectors according to their eigen values in decreasing order.
- Choose first k eigen vectors and that will be the new k dimensions.
- Transform the original n dimensional data points into k dimensions.

Example:

consider the following data

	1	2	3	4	5	6	7	8	9	10
x	2.5	0.5	2.7	1.9	3.1	2.3	-2	1	1.5	1.3
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

Height

$\gamma = 0.926$

Step 1 :- $\bar{x} = 1.81$, $\bar{y} = 1.9$

x	2.5	0.5	2.7	1.9	3.1	2.3	-2	1	1.5	1.3
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

$$X = \begin{bmatrix} x & y \end{bmatrix}_{10 \times 2}$$

Step 2 :-

$$\text{cov}(x) = \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix}$$

$$\text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

$$(A - x\lambda)^T = 0 \Rightarrow \lambda = ? = \frac{1}{n-1} \sum \lambda_i$$

$$Ax = x\lambda$$

Step 3.

Eigen values (λ) = $\lambda_1 = 1.2840$, $\lambda_2 = 0.0490$

Eigen vectors corresponding

$$\lambda_1 \rightarrow \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix} \quad \text{Eigen vector } 1.2840 \rightarrow PC_1$$

$$\lambda_2 \rightarrow \begin{bmatrix} 0.735 \\ -0.678 \end{bmatrix} \quad \text{Eigen vector } 0.0490 \rightarrow PC_2$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 517

i.e.

Variable	Eigen vector 1	Eigen vector 2
x_1	0.678	0.735
x_2	0.735	-0.678

% of total variance

$$\frac{1.2840}{1.333} = 96.3\%$$

$$\frac{0.0490}{1.333} = 3.7\%$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 518

Step 4: select Variation matrix

$$V = \begin{pmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{pmatrix}$$

or

$$V = \begin{pmatrix} 0.678 \\ 0.735 \end{pmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 519

Step 5: Find new data set $Y = XV$

i.e. $Y = X^T V$

Case 1 :- $Y = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{bmatrix}_{2 \times 2}$

$$Y = \begin{bmatrix} 3.459 & 0.211 \\ -0.854 & -0.107 \\ \vdots & \vdots \\ 1.407 & 0.199 \end{bmatrix}_{10 \times 2}$$

$y_1 = 0.678 x_1 + 0.735 x_2$

$y_2 = 0.735 x_1 - 0.678 x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 520

Step 5. Find new Data set $\gamma = xv$

i.e. $\gamma = xv$

case (ii) :- $\gamma = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}$

$= \begin{bmatrix} 3.459 \\ -0.854 \\ \vdots \\ 1.407 \end{bmatrix}$

$\gamma = 0.678x_1 + 0.735x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 521

Summary :-

- 1) Re centre the original data set to the origin
- 2) Find covariance matrix x
- 3) Find eigen values and eigen vectors and also % of variability
- 4) Find the transformation matrix v based on PC selection
- 5) Derive the new data set by $\gamma = xv$

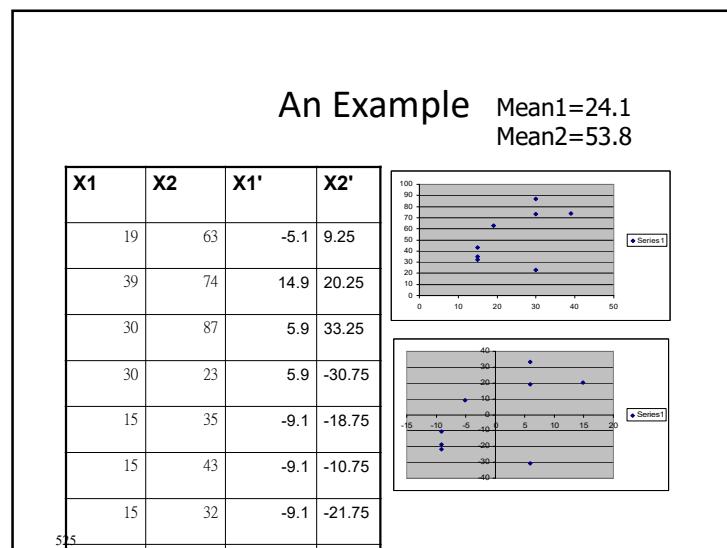
Advanced Statistical Techniques for Analytics 26-04-2019 Slide 522

Principal Components

- All principal components (PCs) start at the origin of the ordinate axes.
- First PC is direction of maximum variance from origin
- Subsequent PCs are orthogonal to 1st PC and describe maximum residual variance

Principal Components

- All principal components (PCs) start at the origin of the ordinate axes.
- First PC is direction of maximum variance from origin
- Subsequent PCs are orthogonal to 1st PC and describe maximum residual variance



The first PC is the linear combination that captures the maximum variance in the data. The second PC is created by selecting another linear combination that max. variance with the constraint that its direction is perpendicular to the first component.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 526



Agenda

- Bi Variate distribution
- multivariate normal distribution
- Eigen values and eigen vectors
- Principal Component Analysis

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 528

Joint Discrete Distribution Function _Discrete

$P(x_i, y_j) \rightarrow$ joint probability distribution function

$P_x(x_i, y_j) \rightarrow$ marginal prob. dist. fun of x

$P_y(x_i, y_j) \rightarrow$ " " " of y .

Joint Discrete Distribution Function _Continuous

$f(x, y) \rightarrow$ joint p.d.f

$f_x(x, y) \rightarrow$ Marginal p.d.f of x

$f_y(x, y) \rightarrow$ Marginal p.d.f of y

Example:-

$$f(x, y) = \begin{cases} Kxy & ; 0 < x < 4, 1 < y < 5 \\ 0 & ; \text{otherwise} \end{cases}$$

1) Find 'K' value?

2) $P(1 < x < 2, 2 < y < 3)$

3) $P(x \geq 3, y \leq 2)$

4) $f_x(x, y)$

5) $f_y(x, y)$

$$1) \int_{x=0}^4 \int_{y=1}^5 f(x, y) dy dx = 1$$

$$\Rightarrow K \int \int xy dy dx = 1$$

$$\Rightarrow K \left[\frac{x^2}{2} \right]_{x=0}^4 \left[\frac{y^2}{2} \right]_{y=1}^5 = 1$$

$$\Rightarrow K \left[\frac{16}{2} - 0 \right] \left[\frac{25}{2} - \frac{1}{2} \right] = 1$$

$$\Rightarrow K = 1/96$$

2) $P(x \geq 3, y \leq 2)$

$$\begin{aligned}
 &= \int_{x=3}^4 \int_{y=1}^2 kxy \, dy \, dx \\
 &= \frac{1}{96} \left[\frac{x^2}{2} \right]_{x=3}^4 \left[\frac{y^2}{2} \right]_{y=1}^2 \\
 &= \frac{1}{96 \times 2 \times 2} \left[4^2 - 3^2 \right] \left[2^2 - 1^2 \right] \\
 &= \frac{7}{128}
 \end{aligned}$$

$$\int x^m dx = \frac{x^{m+1}}{m+1}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 533

3) $P(1 < x < 2, 2 < y < 3)$

$$\begin{aligned}
 &= \int_{x=1}^2 \int_{y=2}^3 kxy \, dx \, dy \\
 &= \frac{1}{96} \left[\frac{x^2}{2} \right]_{x=1}^2 \left[\frac{y^2}{2} \right]_2^3 \\
 &= \frac{5}{128}
 \end{aligned}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 534

4) Marginal of X

$$\begin{aligned}
 f_x(x, y) &= \int_{y=1}^5 f(x, y) \, dy \\
 &= \int_{y=1}^5 kxy \, dy = \frac{1}{96} x \left[\frac{y^2}{2} \right]_1^5 \\
 &= \frac{x}{192} (25 - 1) = \frac{x}{8}
 \end{aligned}$$

$$i.e. f_x(x, y) = \frac{x}{8}, 0 < x < 4$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 535

5) Marginal of Y .

$$\begin{aligned}
 f_y(x, y) &= \int_{x=0}^4 kxy \, dx \\
 &= \frac{1}{96} \left[\frac{x^2}{2} \right]_0^4 \\
 &= \frac{y}{12}
 \end{aligned}$$

$$f_y(x, y) = \frac{y}{12}, 1 < y < 5$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 536

Example:

	0	1	2	
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{2}$	
1	$\frac{2}{9}$	$\frac{1}{6}$	0	
2	$\frac{1}{36}$	0	0	

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 537

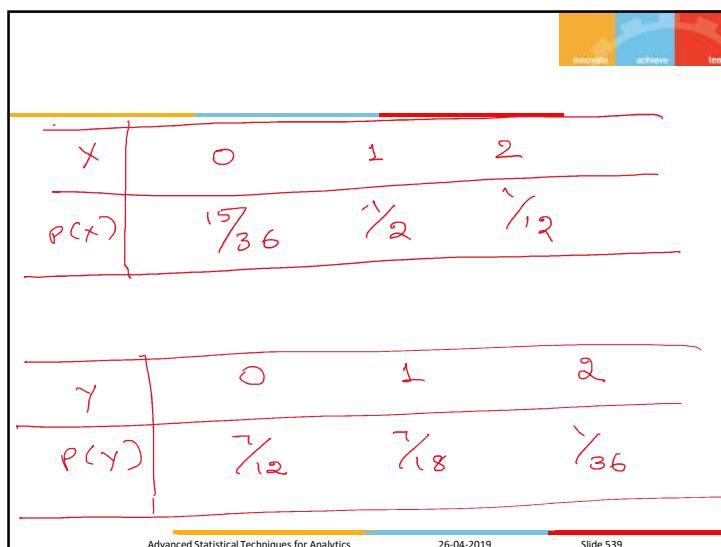
Example:

	0	1	2	Marginal of Y
0	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$
1	$\frac{2}{9}$	$\frac{1}{6}$	0	$\frac{1}{18}$
2	$\frac{1}{36}$	0	0	$\frac{1}{36}$
x	$\frac{15}{36}$	$\frac{1}{2}$	$\frac{1}{12}$	

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 538



Advanced Statistical Techniques for Analytics

26-04-2019

Slide 539

Conditional Probabilities

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(x|y) = \frac{P(x,y)}{P_y(x,y)} \rightarrow \text{joint} \quad \text{Marginal of } Y$$

$$P(y|x) = \frac{P(x,y)}{P_x(x,y)} \rightarrow \text{joint} \quad \text{Marginal of } X$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 540

innovate achieve lead

$$f(x|y) = \frac{f(x,y)}{f_y(x,y)} \xrightarrow{\text{joint}} \text{marginal of } y$$

$$f(y|x) = \frac{f(x,y)}{f_x(x,y)} \xrightarrow{\text{joint}} \text{marginal of } x$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 541

innovate achieve lead

Independent \Rightarrow V.s.

$$P(x,y) = P_x(x,y) \cdot P_y(x,y)$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\text{joint} \quad \text{marginal } x \quad \text{marginal } y$$

$$f(x,y) = f_x(x,y) \cdot f_y(x,y)$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 542

innovate achieve lead

Bi-Variate Normal dist

$$f(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{-\gamma}{2(1-\rho^2)}\right]$$

$$\gamma = \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} - 2\rho \frac{(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2}$$

$$\rho = \text{corr}(x_1, x_2) = \frac{\text{cov}(x_1, x_2)}{\sigma_1\sigma_2}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 543

innovate achieve lead

Multivariate Normal distribution

$$f(x) = \left(\frac{1}{2\pi}\right)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}} \exp\left\{-\frac{1}{2}x^\top \Sigma^{-1} x\right\}$$

x_1, x_2, \dots, x_n

Σ \downarrow
Det of
variance-
covariance
matrix

Σ^{-1} \downarrow
inverse of
variance-
covariance
matrix

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 544

innovate achieve lead

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

covariance between x_1 and x_2 .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 545

innovate achieve lead

Covariance Matrix

- Representing Covariance between dimensions as a matrix
- e.g.

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix}$$

- Diagonal is the variances of x, y and z
- $\text{cov}(x,y) = \text{cov}(y,x)$ hence matrix is symmetrical about the diagonal
- N-dimensional data will result in $n \times n$ covariance matrix

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 546

- innovate achieve lead
- A positive value of covariance indicates both dimensions increase or decrease together
 - A negative value indicates while one increases the other decreases, or vice-versa
 - If covariance is zero: the two dimensions are independent of each other .
- Advanced Statistical Techniques for Analytics 26-04-2019 Slide 547

innovate achieve lead

Transformation matrices

Consider: $A \times X = Y$

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix} = 4 \times \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

Square transformation matrix transforms $(3,2)$ from its original location. Now if we were to take a multiple of $(3,2)$

$$2 \times \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 6 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 6 \\ 4 \end{bmatrix} = \begin{bmatrix} 24 \\ 16 \end{bmatrix} = 4 \times \begin{bmatrix} 6 \\ 4 \end{bmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 548

eigenvalue problem

- The eigenvalue problem is any problem having the following form:
- $A \cdot X = \lambda \cdot X$
- A : $n \times n$ matrix
- X : $n \times 1$ non-zero vector
- λ : scalar
- Any value of λ for which this equation has a solution is called the eigenvalue of A and vector v which corresponds to this value is called the eigenvector of A .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 549

consider the matrix

$$\begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} := 9x + 4y \quad 4x + 3y$$

$(0,0) \rightarrow (0,0)$
 $(1,0) \rightarrow (9,4)$
 $(0,1) \rightarrow (4,3)$
 $(-1,0) \rightarrow (-9,-4)$
 $(0,-1) \rightarrow (-4,-3)$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 550

consider the matrix

$$\begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} := 9x + 4y \quad 4x + 3y$$

$$A \cdot X = \begin{bmatrix} 9 & 4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 9x + 4y \\ 4x + 3y \end{bmatrix}$$

$$= 11 \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

$$= 1 \begin{pmatrix} -1 \\ 2 \end{pmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 551

eigenvalue problem

$$\begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix} = 4 \cdot \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \lambda \cdot \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

Therefore, $(3,2)$ is an eigenvector of the square matrix A and 4 is an eigenvalue of A

eigen value eigen vector

Given matrix A , how can we calculate the eigenvector and eigenvalues for A ?

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 552

innovate achieve lead

$Ax = \lambda x$ eigenvalue eigen vector

$$Ax - \lambda I x = 0$$

$$[A - \lambda I] x = 0$$

Identity matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Homogeneous system of equations

Non trivial soln (non-zero soln)

$2x + 3y = 0$ iff $|A - \lambda I| = 0$

$3x + 2y = 0$

$x_1 = k, x_2 = -k$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 553

innovate achieve lead

Ex: $A = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$

$$|A - \lambda I| = 0 \Rightarrow \begin{vmatrix} 0-\lambda & 1 \\ -2 & -3-\lambda \end{vmatrix} = 0$$

$$\Rightarrow -\lambda(-3-\lambda) + 2 = 0$$

$$\Rightarrow \lambda^2 + 3\lambda + 2 = 0 \quad i.e. \lambda = -1, -2$$

Now we need to find x corresponding to $\lambda = -1$ and $\lambda = -2$ such that $Ax = \lambda x$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 554

innovate achieve lead

when $\lambda = -1$: $[A - \lambda I] x = 0$

$$\therefore \begin{bmatrix} 0+1 & 1 \\ -2 & -2 \end{bmatrix} x = 0 \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$x_1 + x_2 = 0$ and $-2x_1 - 2x_2 = 0$ $\Rightarrow x_1 + x_2 = 0$

$\lambda = -2$

$x_1 = k, x_2 = -k$

$x_1 = \begin{bmatrix} k \\ -k \end{bmatrix}$, $x_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ eigen vectors

Similarly we find eigen vectors corresponding to $\lambda = -2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 555

innovate achieve lead

Principal Component Analysis

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 556

→ why PCA

↓

Dimensionality Reduction

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 557

why PCA ?

→ with large number of variables, the matrix may be too large to study and interpret.

→ Visualization of the data is not helpful and complex also

→ If it is possible to reduce the number of variables to a few, interpretable linear combinations, then it becomes simple.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 558

Applications

- Face Recognition
- Image Compression
- Gene Expression Analysis
- Data Reduction
- Data Classification
- Trend Analysis
- Factor Analysis
- Noise Reduction

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 559

Principal Component Analysis

In real world data analysis tasks we analyze complex data i.e. multi dimensional data. We plot the data and find various patterns in it or use it to train some machine learning models. One way to think about dimensions is that suppose you have a data point x , if we consider this data point as a physical object then dimensions are merely a basis of view, like where is the data located when it is observed from horizontal axis or vertical axis.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 560

The header features a decorative bar divided into three horizontal sections: orange on top, blue in the middle, and red at the bottom. The word "inspire" is written in white on the orange section, "achieve" in white on the blue section, and "lead" in white on the red section.

As the dimensions of data increases, the difficulty to visualize it and perform computations on it also increases. So, how to reduce the dimensions of a data-

- * Remove the redundant dimensions
- * Only keep the most important dimensions

The header bar consists of three horizontal squares. The left square is orange with the word "Dimensionality" in white. The middle square is blue with the word "Analytics" in white. The right square is red with the word "Learning" in white.

- It is a linear transformation that chooses a new coordinate system for the data set such that
 - greatest variance by any projection of the data set comes to lie on the first axis (then called the first principal component),
 - the second greatest variance on the second axis, and so on.
- PCA can be used for reducing dimensionality by eliminating the later principal components.



what does Principal Component Analysis (PCA) do?

PCA finds a new set of dimensions (or a set of basis of views) such that all the dimensions are orthogonal (and hence linearly independent) and ranked according to the variance of data along them. It means more important principle axis occurs first. (more important = more variance/more spread out data)



innovate achieve lead

How does PCA work

- Calculate the covariance matrix X of data points.
- Calculate eigen vectors and corresponding eigen values.
- Sort the eigen vectors according to their eigen values in decreasing order.
- Choose first k eigen vectors and that will be the new k dimensions.
- Transform the original n dimensional data points into k dimensions.

Example:

consider the following data

x	2.5	0.5	2.7	1.9	3.1	2.3	2	1	1.5	1.3
y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

$r = 0.926$

Step 1 :- $\bar{x} = 1.81$, $\bar{y} = 1.9$

x	0.6 0	-1.3 1	0.3 0	0.9 0	0.2 1	0.5 2	0.9 1	0.19 0	-0.88 -1	0.7 0	0.7 -1
y	0.7 9	-1.21 -1	0.9 0	0.29 0	0.9 1	0.79 0	0.19 1	0.81 0	0.21 -1	0.1 0	0.1 -1

$X = \begin{bmatrix} x & y \end{bmatrix}$
 10×2

Step 2 :-

$$\text{cov}(x) = \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix}$$

$$\text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

$$|A - \lambda I| = 0 \Rightarrow \lambda = ?$$

$$\lambda = \frac{1}{n-1} \sum \lambda_i$$

$$AX = \lambda x$$

Step 3 :-

Eigen values (x) = $\lambda_1 = 1.2840$, $\lambda_2 = 0.0490$

Eigen vectors corresponding

$$\lambda_1 \rightarrow \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix} \quad \checkmark 1.2840 \rightarrow PC_1$$

$$\lambda_2 \rightarrow \begin{bmatrix} 0.735 \\ -0.678 \end{bmatrix} \quad \checkmark 0.0490 \rightarrow PC_2$$

ie

Variable	Eigen vector ₁ $\lambda_1 = 1.2840$	Eigen vector ₂ $\lambda_2 = 0.0490$	Total 1.333
x_1	0.678	0.735	
x_2	0.735	-0.678	
% of Total Variance	1.2840 = 96.3%	0.0490 = 3.7%	

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 569

Step 4: select Variation matrix

$$V = \begin{pmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{pmatrix}$$

or

$$V = \begin{pmatrix} 0.678 \\ 0.735 \end{pmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 570

Step 5. Find new data set $Y = XV$

ie $Y = XV$

case(i) :- $Y = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{bmatrix}_{2 \times 2}$

$\therefore Y = \begin{bmatrix} 3.459 & 0.211 \\ -0.854 & -0.107 \\ \vdots & \vdots \\ 1.407 & 0.199 \end{bmatrix}_{10 \times 2}$

$y_2 = 0.735x_1 - 0.678x_2$

$y_1 = 0.678x_1 + 0.735x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 571

Step 5. Find new data set $Y = XV$

ie $Y = XV$

case(ii) :- $Y = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}_{2 \times 1}$

$\therefore Y = \begin{bmatrix} 3.459 \\ -0.854 \\ \vdots \\ 1.407 \end{bmatrix}_{10 \times 1}$

$\therefore Y = 0.678x_1 + 0.735x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 572

Summary :-

- 1) Re centre the original data set to the origin
- 2) Find covariance matrix Σ
- 3) Find eigen values and eigen vectors and also % of variability
- 4) Find the transformation matrix V based on PC selection
- 5) Derive the new data set by $Y = \Sigma V$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 573

Principal Components

- All principal components (PCs) start at the origin of the ordinate axes.
- First PC is direction of maximum variance from origin
- Subsequent PCs are orthogonal to 1st PC and describe maximum residual variance

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 573

Principal Components

- All principal components (PCs) start at the origin of the ordinate axes.
- First PC is direction of maximum variance from origin
- Subsequent PCs are orthogonal to 1st PC and describe maximum residual variance

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 573

The first PC is the linear combination that captures the maximum variance in the data. The second PC is created by selecting another linear combination that max. variance with the constraint that its direction is perpendicular to the first component.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 576

Question :

How many principal Components should be considered?

Any criteria ---?

Component	Eigen value	Proportion
1	0.3775	0.7227%
2	0.0511	0.0977 82%
3	0.0279	0.0535 87%
4	0.0230	0.0440 → 91
5	0.0168	0.0321 → 95
6	0.0120	0.0229 → 97
7	0.0085	0.0162 → 98
8	0.0039	0.0075 → 99
9	0.0018	0.0034 → 100
Total		0.5225

1) First two eigen values together explains 82% of the variation.

First three eigen values explains 87% of the variation

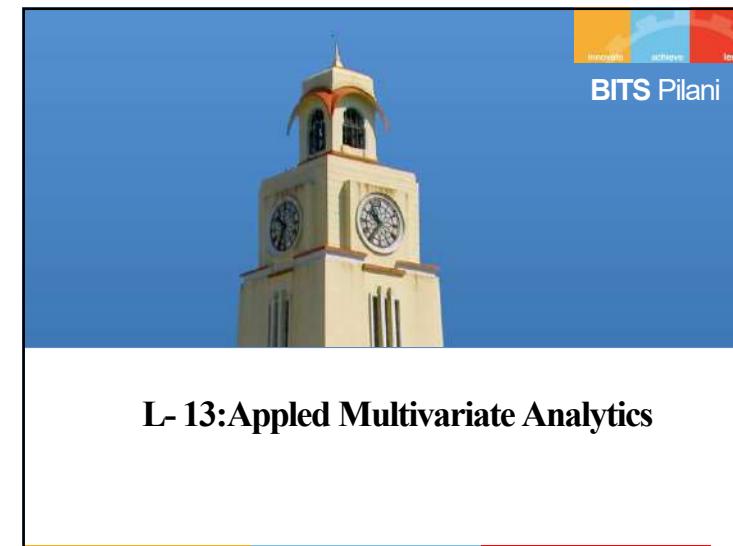
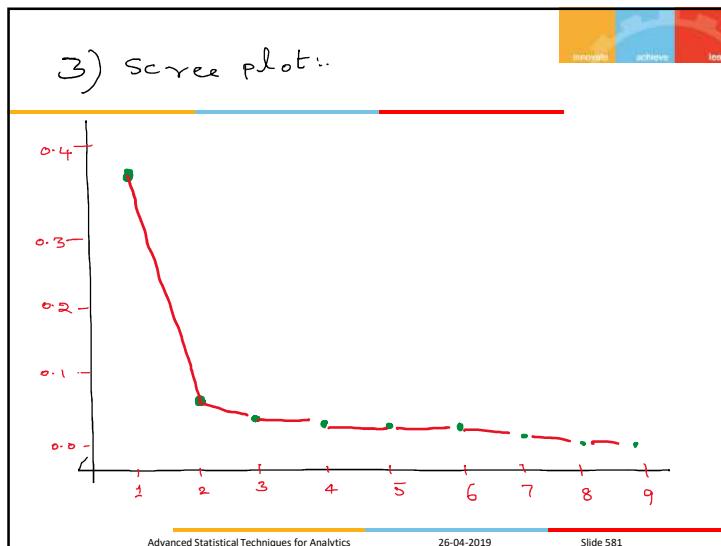
$\lambda_1 \rightarrow PC_1$
 $\lambda_2 \rightarrow PC_2$

$\lambda_1 \rightarrow PC_1$
 $\lambda_2 \rightarrow PC_2$
 $\lambda_3 \rightarrow PC_3$

2) $\lambda_1 = 0.3775 \}$ → 0.326
 $\lambda_2 = 0.0511 \}$ → 0.0232
 $\lambda_3 = 0.0279 \}$ → 0.0049
 $\lambda_4 = 0.0230 \}$ → 0.0049

Like these, the differences are smaller and smaller.

This is another indicator of how many eigen values to consider



Agenda_13th Session

- Principal Component Analysis
- ANOVA

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 583

Principal Component Analysis

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 584

Example:

consider the following data

<u>x</u>	2.5	0.5	2.7	1.9	3.1	2.3	2	1	1.5	1.3
<u>y</u>	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9

$$r = 0.926$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 585

Research | Achieve | Lead

Step 1 :- $\bar{x} = 1.81$, $\bar{y} = 1.91$

x	y	x_1	y_1	x_2	y_2	x_3	y_3	x_4	y_4	x_5	y_5
0.69	0	-0.21	-1.31	0.39	0.09	0.29	0	1.09	1.29	0.79	0.19
0.49	-0.19	0.99	-1.21	-0.29	0.39	-0.29	0.99	-0.99	-0.29	-0.31	-0.81
0.29	-0.31	-0.99	0.69	0.99	-0.29	-0.99	-0.99	0.29	-0.99	-0.31	-0.31
0.09	-0.51	-1.69	0.49	-0.99	0.29	-0.99	-0.99	0.99	-0.99	-0.19	-0.71
-0.09	-0.71	0.89	-0.49	0.29	-0.29	0.29	-0.99	0.99	-0.99	-0.31	-0.19
-0.29	-0.31	-0.99	0.69	-0.99	0.29	-0.99	-0.99	0.99	-0.99	-0.31	-0.31
-0.49	-0.51	0.29	-0.49	0.99	-0.29	-0.99	-0.99	0.99	-0.99	-0.19	-0.71
-0.69	-0.71	-0.99	0.49	-0.99	0.29	-0.99	-0.99	0.99	-0.99	-0.31	-0.31
-0.89	-0.91	-1.69	0.29	-0.99	0.29	-0.99	-0.99	0.99	-0.99	-0.19	-0.71
-1.09	-1.11	-1.89	0.09	-0.99	0.29	-0.99	-0.99	0.99	-0.99	-0.31	-0.31

$$X = \begin{bmatrix} x & y \\ \end{bmatrix}$$

10×2

Integrity achieve

Step 2 :-

$$\text{cov}(x) = \begin{bmatrix} 0.6166 & 0.6154 \\ 0.6154 & 0.7166 \end{bmatrix}$$

↓

$$\text{cov}(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

$$|A - x\bar{x}| = 0 \rightarrow x = ? : \frac{1}{n-1} \sum xy$$

\checkmark $\boxed{Ax = x\bar{x}}$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 587

Step 3:

$$\text{Eigen values } (\lambda) = \lambda_1 = 1.2840 \\ \lambda_2 = 0.0490$$

Eigen vectors corresponding

$$\lambda_1 \rightarrow$$

$$\begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}$$

$$1.2840$$

PC₁

$$\lambda_2 \rightarrow$$

$$\begin{bmatrix} 0.735 \\ -0.678 \end{bmatrix}$$

$$0.0490$$

PC₂

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 588

ie

Variable	Eigen vector ₁ $\lambda_1 = 1.2840$	Eigen vector ₂ $\lambda_2 = 0.0490$	Total 1.333
x_1	0.678	0.735	
x_2	0.735	-0.678	
% of Total Variance	1.2840	0.0490	
	1.333	0.333	
	= 96.3%	= 3.7%	

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 589

Step 4: select Variation matrix

$$V = \begin{pmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{pmatrix}$$

or

$$V = \begin{pmatrix} 0.678 \\ 0.735 \end{pmatrix}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 590

Step 5. Find new data set $Y = XV$

ie $Y = XV$

case(i) :- $Y = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 & 0.735 \\ 0.735 & -0.678 \end{bmatrix}_{2 \times 2}$

$\therefore Y = \begin{bmatrix} 3.459 & 0.211 \\ -0.854 & -0.107 \\ \vdots & \vdots \\ 1.407 & 0.199 \end{bmatrix}_{10 \times 2}$

$y_2 = 0.735x_1 - 0.678x_2$

$y_1 = 0.678x_1 + 0.735x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 591

Step 5. Find new data set $Y = XV$

ie $Y = XV$

case(ii) :- $Y = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ \vdots & \vdots \\ 1.1 & 0.9 \end{bmatrix}_{10 \times 2} \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}_{2 \times 1}$

$\therefore Y = \begin{bmatrix} 3.459 \\ -0.854 \\ \vdots \\ 1.407 \end{bmatrix}_{10 \times 1}$

$\therefore Y = 0.678x_1 + 0.735x_2$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 592

Summary :-

- 1) Re centre the original data set to the origin
- 2) Find covariance matrix Σ
- 3) Find eigen values and eigen vectors and also % of variability
- 4) Find the transformation matrix V based on PC selection
- 5) Derive the new data set by $Y = XU$

The first PC is the linear combination that captures the maximum variance in the data. The second PC is created by selecting another linear combination that max. variance with the constraint that its direction is perpendicular to the first component.

Question :

How many principal Components should be considered?

Any criteria ---?

Component	Eigen value	Proportion
1	0.3775	0.7227%
2	0.0511	0.0977% 82%
3	0.0279	0.0535% 87%
4	0.0230	0.0440% 91
5	0.0168	0.0321% 95
6	0.0120	0.0229% 97
7	0.0085	0.0162% 98
8	0.0039	0.0075% 99
9	0.0018	0.0034% 100
Total		0.5225

1) First two eigen values together explains 82% of the variation.
 First three eigen values explain 87% of the variation.

$\lambda_1 \rightarrow PC_1$

$\lambda_2 \rightarrow PC_2$

$\lambda_1 \rightarrow PC_1$

$\lambda_2 \rightarrow PC_2$

$\lambda_3 \rightarrow PC_3$

PC	Eigenvalue
1	0.3775
2	0.0511
3	0.0279
4	0.0230
5	0.0150
6	0.0120
7	0.0090
8	0.0070
9	0.0050
10	0.0030

2) $\lambda_1 = 0.3775$

$\lambda_2 = 0.0511$

$\lambda_3 = 0.0279$

$\lambda_4 = 0.0230$

$\lambda_1 = 0.3775 \rightarrow 0.326$

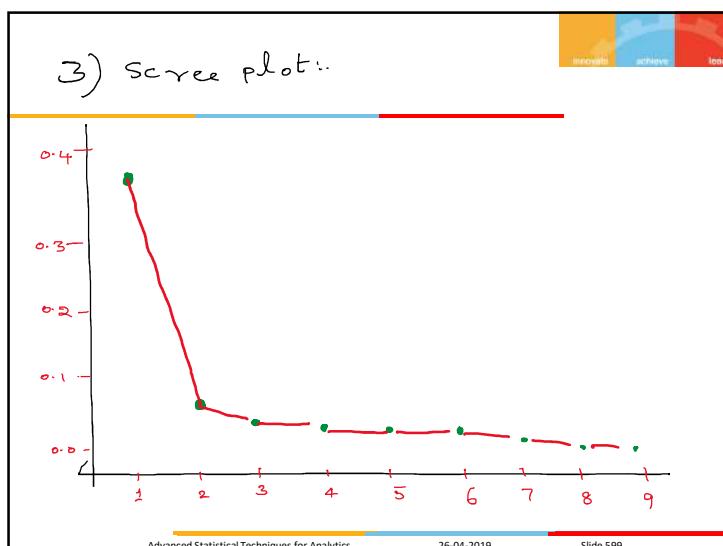
$\lambda_2 = 0.0511 \rightarrow 0.0232$

$\lambda_3 = 0.0279 \rightarrow 0.0049$

$\lambda_4 = 0.0230 \rightarrow 0.0049$

Like these, the differences are smaller and smaller.

This is another indicator of how many eigen values to consider



ANOVA

- ✓ Effectiveness of different promotional activities
- ✓ Quality of a product produced by different manufacturers in terms of an attribute
- ✓ Yield of crop due to varieties of seeds , fertilisers and quality of soil

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 601

Between/Within Groups

Variance can be separated into two major components

- **Within groups** – variability or differences in particular groups (individual differences)
- **Between groups** - differences depending what group one is in or what treatment is received

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 602

Null hypothesis:
no difference in means

Alternative hypothesis:
difference in means

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 603

ANOVA

Effectiveness of different promotional activities

Quality of a product produced by different manufacturers in terms of an attribute

Yield of crop due to varieties of seeds , fertilisers and quality of soil

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 604

ANOVA-analysis of variance

* Significance of difference between two sample means

$$H_0 = \mu_1 = \mu_2 = \dots = \mu_K$$

$$H_1 = \mu_1 \neq \mu_2 \neq \dots \neq \mu_K$$

Alternative Hypothesis

Assumptions

Each population is normally distributed with mean μ_i With equal variances

$$\sigma^2_i$$

Each sample is drawn randomly and independent of other samples

Observations	population		
	A	B	C
1	26	18	23
2	25	16	19
3	28	17	26
4	12	8	30

ANOVA summary

Source of Variation	Sum of squares	d.o.f	Mean squares	F-value
Between Samples	SSTR	n-1	MSTR = SSTR / (n-1)	
within Samples (Error)	SSE	n-n	MSE = SSE / (n-n)	$F = \frac{MSTR}{MSE}$
Total	SST	n-1		

Short cut method

$$T = \sum x_1 + \sum x_2 + \dots + \sum x_n$$

Corr. Fact = CF = $\frac{T^2}{m}$; $m = m_1 + m_2 + \dots + m_n$

$$SST = \left[\sum (x_1^2) + \sum (x_2^2) + \dots + \sum (x_n^2) \right] - CF$$

$$SSTR = \frac{\left(\sum x_j \right)^2}{m} - CF$$

$$SSE = SST - SSTR$$

Example

A study was conducted to investigate the perception of corporate ethical values among individuals specialising in marketing. Using 0.05 level of significance and the data given below, test for significant differences in perception among three groups. (higher scores indicate higher ethical values)

	Marketing manager	Marketing Research	Advertising
1	6	5	6
2	5	5	7
3	4	4	6
4	5	4	5
5	6	5	6
6	4	4	6

$$n=3, m=18$$

$$T = \sum x_1 + \sum x_2 + \sum x_3 \\ = 30 + 27 + 36 = 93$$

$$CF = \frac{T^2}{m} = \frac{(93)^2}{18} = 480.50$$

$$SST = (\sum x_1^2 + \sum x_2^2 + \sum x_3^2) - CF \\ = 154 + 123 + 218 = 495.50$$

innovate achieve lead

$$\begin{aligned}
 SSTR &= \left(\frac{\sum x_1^2}{m_1} + \frac{\sum x_2^2}{m_2} + \frac{\sum x_3^2}{m_3} \right) - CF \\
 &= \frac{(30)^2}{6} + \frac{(27)^2}{6} + \frac{(36)^2}{6} - 480.50 \\
 &\approx 7 \\
 \checkmark SSE &= SST - SSTR \\
 &= 14.50 - 7 = 7.50
 \end{aligned}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 613

innovate achieve lead

$$\begin{aligned}
 MSTR &= \frac{SSTR}{df_1} = \frac{7}{2} : 3.5 \\
 MSE &= \frac{SSE}{df_2} = \frac{7.5}{df_2} : 0.5 \\
 F &= \frac{MSTR}{MSE} = \frac{3.5}{0.5} = 7 \quad (\text{F-distribution table}) \\
 \text{calculated value: } 7 & \\
 \text{F-table value: } 3.68 & \\
 3.68 < 7 \Rightarrow & \text{Rejected}
 \end{aligned}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 614

Two way ANOVA

Sources of variation	Sum of square	D.F.	mean square	test statistic
Between columns }	SSTR	C-1	MSTR = $\frac{SSTR}{C-1}$	$F_{\text{Treatment}}$
Between rows }	SSR	n-1	MSR = $\frac{SSR}{n-1}$	$F_{\text{MSR/MSE}}$
Residual error }	SSE	(C-1)(n-1)	MSE = $\frac{SSE}{(C-1)(n-1)}$	F_{Blocks}
Total	SST	n-1		$F_{\text{MSR/MSE}}$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 615

Example

Two way ANOVA

Month	A	B	C	D	Sales
May	50	40	48	39	
June	46	48	50	45	
July	39	44	40	39	

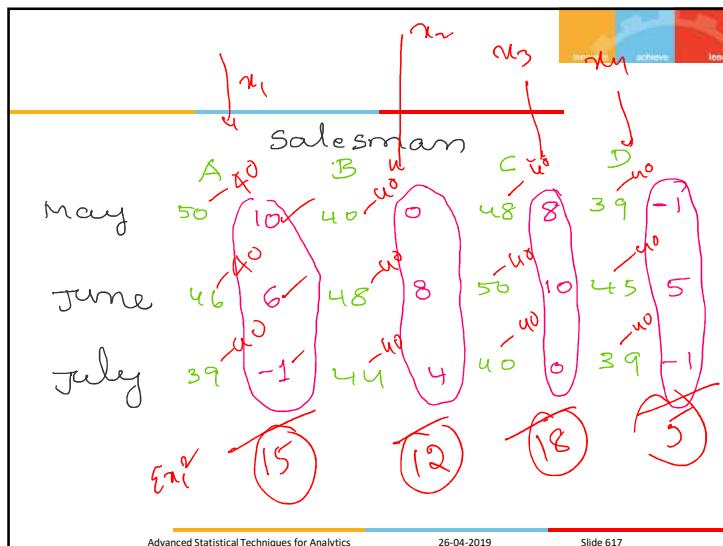
→ Is there any significant difference in the sales made, if so,

→ Is there a significant diff in the sales made during these months.

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 616



$$\begin{aligned}
 T &= 15 + 12 + 18 + 5 = 48 \\
 CF &= \frac{T^2}{n} = \frac{(48)^2}{12} = 192 \\
 SSTR &= \text{Sum of squares (columns)} \\
 &= \left(\frac{15^2}{3} + \frac{12^2}{3} + \frac{18^2}{3} + \frac{5^2}{3} \right) - 192 \\
 &= 42
 \end{aligned}$$

$$\begin{aligned}
 SSR &= \text{Sum of squares between months (rows)} \\
 &= \left(\frac{17^2}{4} + \frac{29^2}{4} + \frac{22^2}{4} \right) - 192 \\
 &= 91.5 \\
 SST &= \left(\sum x_1^2 + \sum x_2^2 + \sum x_3^2 + \sum x_4^2 \right) - CF \\
 &= (137 + 80 + 164 + 27) - 192 \\
 &= 216
 \end{aligned}$$

$$\begin{aligned}
 SSE &= SST - (SSTR + SSR) \\
 &= 216 - (42 + 91.5) \\
 &= 82 \\
 df_c &: 3, \quad df_{n-1} : n-1 : 2 \\
 df &: (c-1)(n-1) : 3 \times 2 = 6
 \end{aligned}$$

innovate achieve lead

$$MSTR = \frac{SSTR}{c-1} = \frac{42}{3} = 14$$

$$MSR = \frac{SSR}{n-1} = \frac{91.5}{12} = 45.75$$

$$MSE = \frac{SSE}{(c-1)(n-1)} = \frac{82.5}{6} = 13.75$$

s

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 621

innovate achieve lead

	MSTR < MSE	F = MSE / MSTR	MSTR > MSE	F = MSTR / MSE	Welchman
					Variance ratio
					F treatment
x Between salesmen	SSTR 42.0 ✓	c-1 3	MSTR $\frac{42.0}{3} = 14.$		$F = \frac{14}{13.75} = 1.018$
x Between months	SSR 91.5	n-1 2	MSR = 45.75		$F_{block} = \frac{45.75}{13.75} = 3.327$
x residual error	SSE 82.5	(c-1)(n-1) 6	MSE 13.75		row ² months
Total	216	11			

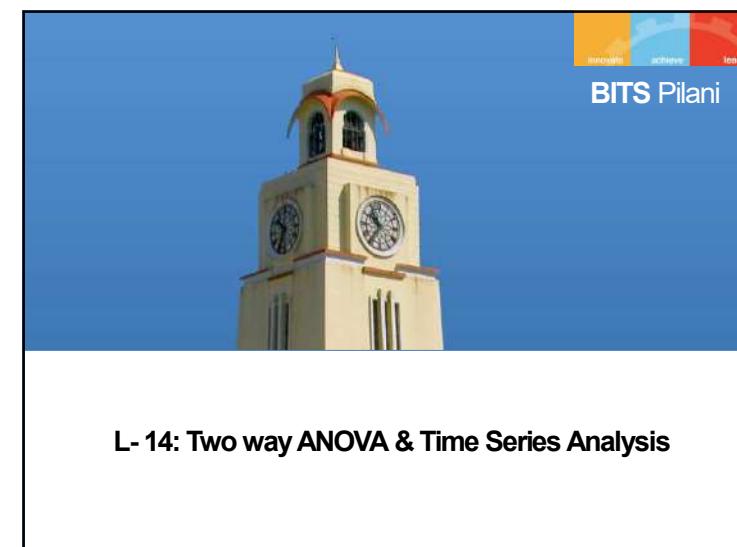
Advanced Statistical Techniques for Analytics 26-04-2019 Slide 622

innovate achieve lead

- $F_{treatment} = 1.018 < F_{0.05}$ (df₁:3, df₂:6) accept
- $F_{block} = 3.327 < F_{0.05}$ (2, 6) accept

difference in the sales by salesman
difference in sales made during months.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 623



Two way ANOVA			
Sources of variation	Sum of square	D.F	mean square
Between columns }	SSTR	C-1	MSTR = $\frac{SSTR}{C-1}$
Between rows }	SSR	n-1	MSR = $\frac{SSR}{n-1}$
Residual error }	SSE	(C-1)(n-1)	MSE = $\frac{SSE}{(C-1)(n-1)}$
Total	SST	n-1	

test statistic

F-treatment
F-blocks

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 625

Example				
Month	A	B	C	
May	50	40	48	39
June	46	48	50	45
July	39	44	40	39

→ IS there any significant difference in the sales made.

✓ IS there a significant diff in the sales made during these months.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 626

Salesman			
May	A 40	B 40	C 48
June	46	48	50
July	39	44	40
	15	12	18

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 627

$$\begin{aligned}
 T &= 15 + 12 + 18 + 3 = 48 \\
 CF &= \frac{T^2}{n} = \frac{48^2}{12} = 192 \\
 SSTR &= \text{Sum of squares (columns)} \\
 &= \left(\frac{15^2}{3} + \frac{12^2}{3} + \frac{18^2}{3} + \frac{3^2}{3} \right) - 192 \\
 &= 42
 \end{aligned}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 628

SSR = Sum of squares between months (rows)

$$= \left(\frac{17^2}{4} + \frac{29^2}{4} + \frac{22^2}{4} \right) - 192$$

$$= 91.5$$

$$SST = (\sum x_1^2 + \sum x_2^2 + \sum x_3^2 + \sum x_4^2) - CF$$

$$= (137 + 80 + 164 + 27) - 192$$

$$= 216$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 629

$$\begin{aligned} SSE &= SST - (SSTR + SSR) \\ &= 216 - (42 + 91.5) \\ &= 82 \end{aligned}$$

$$df_c : 3, \quad df_{n-1} : n-1 : 3-1 : 2$$

$$df : (c-1)(n-1) : 3 \times 2 = 6$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 630

$$MSTR = \frac{SSTR}{c-1} = \frac{42}{3} = 14$$

$$MSR = \frac{SSR}{n-1} = \frac{91.5}{2} = 45.75$$

$$MSE = \frac{SSE}{(c-1)(n-1)} = \frac{82.5}{6} = 13.75$$

Advanced Statistical Techniques for Analytics

26-04-2019

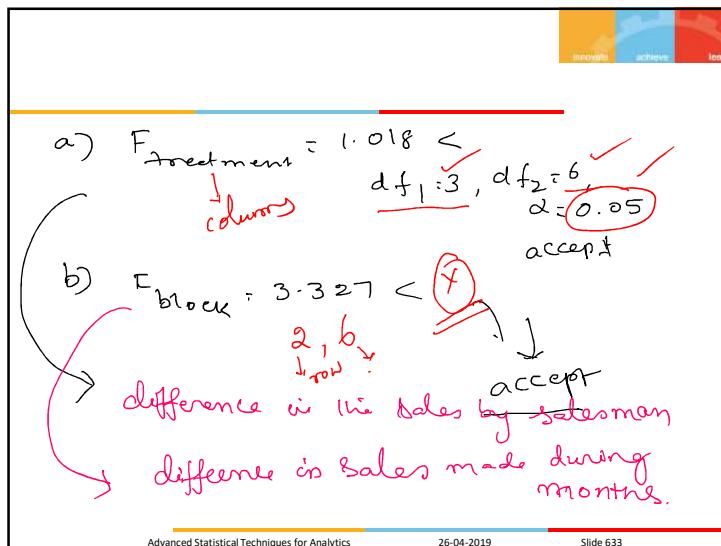
Slide 631

	Sum of squares	D.o.f	mean squares	
x Between Salesmen	SSTR 42.0	c-1 3	MSTR $\frac{42.0}{3} = 14$	F-test $\frac{14}{13.75} = 1.018$
x Between months	SSR 91.5	n-1 2	MSR 45.75	F-block $\frac{45.75}{13.75} = 3.327$
x residual errors	SSE 82.5	(c-1)(n-1) 6	MSE 13.75	rowP columns
Total	216	11		Variance ratio

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 632



Time Series Analysis

- What is time series
- Why time series
- Components in time series
- Methods / models in time series

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 634

What is a Time Series?

Set of evenly spaced numerical data

- Obtained by observing response variable at regular time periods

Forecast based only on past values

- Assumes that factors influencing past, present, & future will continue

Example

Year:	1995	1996	1997	1998	1999
Sales:	78.7	63.5	89.7	93.2	92.1

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 635

Why Time Series?

- To detect patterns of change over regular intervals of time.
- Predict for the future intervals of time

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 636

Applications

- Retail sales
- Spare parts planning
- Stock trading

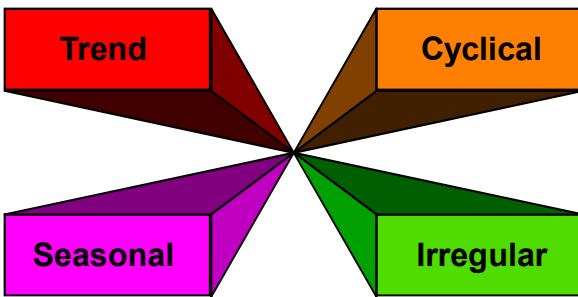
Advanced Statistical Techniques for Analytics 26-04-2019 Slide 637

Time series _ components

- Trend
- Cyclical
- Seasonality
- Random / irregularity

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 638

Time Series Components

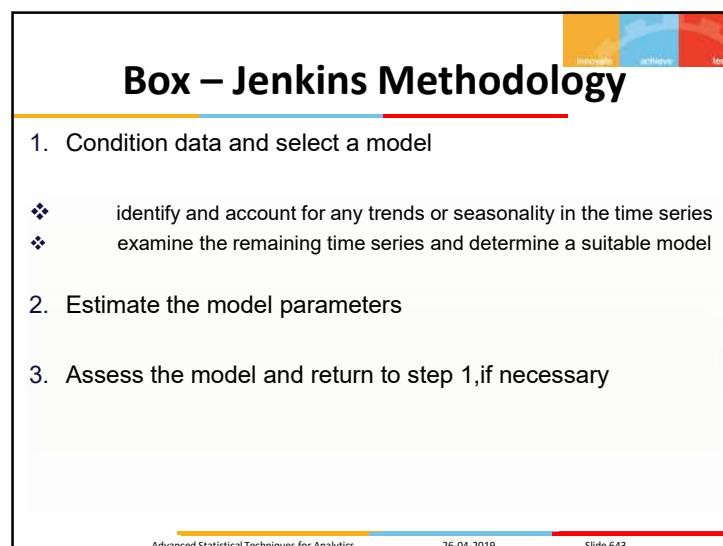
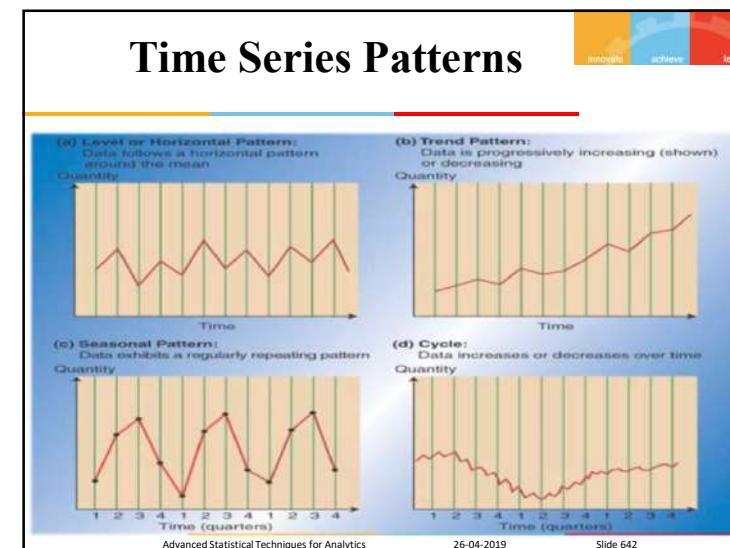
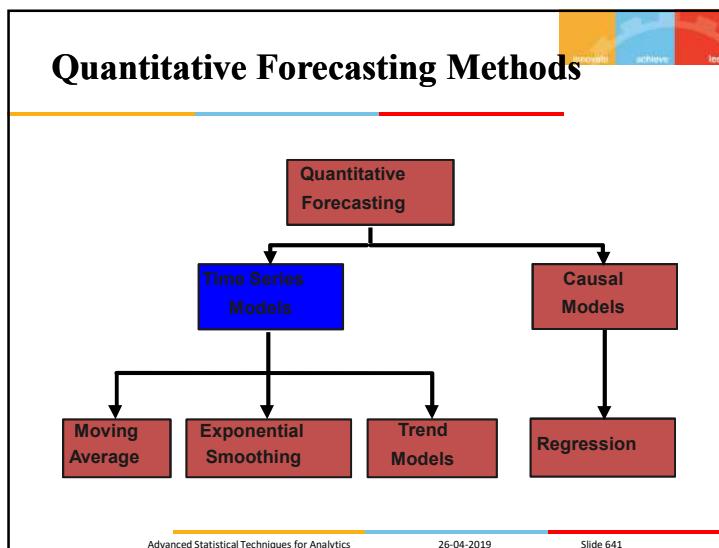


Advanced Statistical Techniques for Analytics 26-04-2019 Slide 639

Time Series Models

- Forecaster looks for data patterns as
Data = historic pattern + random variation
- Historic pattern to be forecasted:
 - Level (long-term average) – data fluctuates around a constant mean
 - Trend – data exhibits an increasing or decreasing pattern
 - Seasonality – any pattern that regularly repeats itself and is of a constant length
 - Cycle – patterns created by economic fluctuations
- Random Variation cannot be predicted

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 640



Time Series Components

A time series can be described by models based on the following components

T_t	Trend Component
S_t	Seasonal Component
C_t	Cyclical Component
I_t	Irregular Component

Using these components we can define a time series as the sum of its components or an **additive model**

$$X_t = T_t + S_t + C_t + I_t$$

Alternatively, in other circumstances we might define a time series as the product of its components or a **multiplicative model** – often represented as a logarithmic model

$$X_t = T_t S_t C_t I_t$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 644

innovate achieve lead

Smoothing Methods

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 645

innovate achieve lead

Moving Average Models

Simple Moving Average Forecast

$$F_t = E(Y_t) = \frac{\sum_{i=t-k}^{t-l} Y_i}{k}$$

Weighted Moving Average Forecast

$$F_t = E(Y_t) = \frac{\sum_{i=t-k}^{t-l} w_i Y_i}{k}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 646

innovate achieve lead

Example(Moving averages)

Use the following data to compute three year moving average for all available years. Find the trend and Forecast error

YEAR	Saleson (Lakhs)	YEAR	Saleson (Lakhs)
2008	21	2013	22
2009	22	2014	25
2010	23	2015	26
2011	25	2016	27
2012	24	2017	26

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 647

innovate achieve lead

Year	Product	3 Year Moving Avg	Error forecast
2008	21	$\overline{21}$	
2009	22	$\overline{21+22} = 22.00$	
2010	23	$\overline{22+23} = 22.50$	-0.33
2011	25	$\overline{23+25} = 24.00$	1.00
2012	24	$\overline{25+24} = 24.50$	-0.33
2013	22	$\overline{24+22} = 23.00$	-1.67
2014	25	$\overline{22+25} = 23.50$	0.67
2015	26	$\overline{25+26} = 25.50$	0.00
2016	27	$\overline{26+27} = 26.50$	-0.67
2017	26		

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 648

Time Series Models

- Weighted Moving Average:**

 - All weights must add to 100% or 1.00
e.g. $C_{t-5} .5, C_{t-1} .3, C_{t-2} .2$ (weights add to 1.0)
 - $F_{t+1} = \sum C_i A_i$
 - Allows emphasizing one period over others; above indicates more weight on recent data ($C_t=0.5$)
 - Differs from the simple moving average that weighs all periods equally - more responsive to trends

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 649

Example(Weighted moving Averages)

Weights	Month
3	Last month
2	Two months ago
1	Three months ago

Months	1	2	3	4	5	6	7	8	9	10	11	12
Sales	10	12	13	16	19	23	26	30	28	18	16	14

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 650

Weights	Example(Weighted moving Averages)
3	Last month
2	Two months ago
1	Three months ago

Months	1	2	3	4	5	6	7	8	9	10	11	12
Sales	10	12	13	16	19	23	26	30	28	18	16	14

Handwritten notes: A 3-month weighted moving average is calculated for month 12. The weights are 1, 2, and 3. The values are 10, 12, and 13 respectively. The calculation is shown as $(10 \times 1) + (12 \times 2) + (13 \times 3) = 12.1$.

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 651

Example.

month	Demand
43	105
44	106
45	110
46	110
47	114
48	121
49	130
50	128
51	137

a) forecast demand for month 52 using 5-months moving Avg

b) " weighted moving average with weights 3, 2, 1 - latents descending

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 652

<u>Example:-</u>		innovate	achieve	lead
month	Demand			
43	105	a)		
44	106	-	114	$f_{t+1} = 121 + 130 \times \frac{1}{5}$
45	110	sus		$128 + 137 \times \frac{1}{5}$
46	110	→ 109.50		= 126
47	119	→ 117.0	b)	$3 + 137 + 2 + 128 + 1 \times 130 \times \frac{1}{6}$
48	121	→ 120.6		= 133 units
49	130	→ 126.0		
50	128	-		
51	137	-		

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 653

Time Series Models

- Exponential Smoothing:

Most frequently used time series method because of ease of use and minimal amount of data needed

- Need just three pieces of data to start:
 - Last period's forecast (F_t)
 - Last period's actual value (A_t)
 - Select value of smoothing coefficient, α , between 0 and 1.0
$$F_{t+1} = \alpha A_t + (1 - \alpha) F_t$$

$$= F_t + \alpha (A_t - F_t)$$
- If no last period forecast is available, average the last few periods or use naive method
- Higher values may place too much weight on last period's random variation

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 654

<u>Example:-</u>		innovate	achieve	lead
Forecast for the first week of March was 500 units whereas the actual demand is 450 units				
a)	Forecast demand for the next week in March			
b)	Assume the actual demand during the March 8 is 505 units.			
	continue the forecasting, assuming that subsequent demands were actually 516, 488, 467, 554 and 510 units.			

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 655

Example:-

Forecast for the first week of March was 500 units whereas the actual demand is 450 units

a) Forecast demand for the next week in March

$$F_{t+1} = F_t + \alpha (A_t - F_t)$$

$$= 500 + 0.1 (450 - 500)$$

$$= 495$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 656

Week	Demand	(A _t) Forecast	(F _t) (old)	new forecast
March 1	450	500	500 + 0.1 (450-500) = 495	
8	505	495	495 + 0.1 (505-495) = 496	
15	516	496	496 + 0.1 (516-496) = 498	
22	488	498	498 + 0.1 (488-498) = 497	
April 1	467	497	497 + 0.1 (467-497) = 494	
8	554	494	494 + 0.1 (554-494) = 500	
15	510	500	500 + 0.1 (510-500) = 501	

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 657

Forecasting Trend

- Basic forecasting models for trends compensate for the lagging that would otherwise occur

- One model, **trend-adjusted exponential smoothing** uses a three step process

- Step 1 - Smoothing the level of the series

$$S_t = \alpha A_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

- Step 2 - Smoothing the trend

$$T_t = \beta(S_t - S_{t-1}) + (1-\beta)T_{t-1}$$

- Forecast including the trend

$$FIT_{t+1} = S_t + T_t$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 658

Measuring Forecasting Accuracy

Mean Absolute Deviation (MAD)

- measures the total error in a forecast without regard to sign

$$MAD = \frac{\sum |actual - forecast|}{n}$$

Cumulative Forecast Error (CFE)

- Measures any bias in the forecast

$$CFE = \sum (actual - forecast)$$

Mean Square Error (MSE)

- Penalizes larger errors

$$MSE = \frac{\sum (actual - forecast)^2}{n}$$

Tracking Signal

- Measures if your model is working

$$TS = \frac{CFE}{MAD}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 659

Stationarity

stationary time series have no trend.

conditions

- constant mean
- Constant variance
- An autocovariance that does not depend on time

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 660

Auto Correlation

auto covariance _{h} (x_t)
 $= \text{cov}(x_t, x_{t-h})$

auto correlation _{h} (x_t)
 $= \frac{\text{Auto cov}_h(x_t)}{\text{std}(x_t) \text{ std}(x_{t-h})}$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 661

Auto Correlation Function

auto covariance _{h} (x_t)
 $\hat{\gamma}_x(h) = \text{cov}(x_t, x_{t-h})$

$\text{ACF} = \frac{\hat{\gamma}_x(h)}{\hat{\gamma}_x(0)} = \text{Cor}(x_t, x_{t-h})$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 662

Models

- AR Model $\rightarrow AR(p)$
- MA Model $\rightarrow MA(q)$
- ARMA Model $\rightarrow ARMA(p,q)$
- ARIMA Model

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 663

AR Model(Auto regressive model)

AR(p)

$$y_t = \delta_t + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \epsilon_t$$

\downarrow Constant \downarrow value of time series at time $t-j$ $\phi_p \neq 0$ $\epsilon_t \sim N(0, \sigma^2)$ for all t .

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 664

Moving Average(MA) Model

$$y_t = f(\epsilon_t, (\epsilon_{t-1}, \epsilon_{t-2}, \dots))$$

today's announcement

$$\text{MA}(\theta) = \theta_0 + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

θ_k is constant for $k=1, 2, \dots, q$

$\theta_q \neq 0$

$\epsilon_t \sim N(0, \sigma^2)$ for all t .

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 665

ARMA model – ARMA(p,q)

$$y_t = \delta + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

If $p=0$ and $q>0$, then AR(p)

If $p=0$ and $q\neq 0$, then MA(q)

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 666

Selecting the Right Forecasting Model

1. The amount & type of available data
 - Some methods require more data than others
2. Degree of accuracy required
 - Increasing accuracy means more data
3. Length of forecast horizon
 - Different models for 3 month vs. 10 years
4. Presence of data patterns
 - Lagging will occur when a forecasting model meant for a level pattern is applied with a trend

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 667

Example:-

Forecast for the first week of March was 500 units whereas the actual demand is 450 units

- a) Forecast demand for the next week i.e. March 8

$$F_{t+1} = F_t + \alpha (A_t - F_t)$$

$$= 500 + 0.1 (450 - 500)$$

$$= 495$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 668

Week	Demand	(A _t) Forecast	(F _t) (old)	new forecast
March 1	450	500	500 + 0.1 (450-500) = 495	
8	505	495	495 + 0.1 (505-495) = 496	
15	516	496	496 + 0.1 (516-496) = 498	
22	488	498	498 + 0.1 (488-498) = 497	
April 1	467	497	497 + 0.1 (467-497) = 494	
8	554	494	494 + 0.1 (554-494) = 500	
15	510	500	500 + 0.1 (510-500) = 501	

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 669

Forecasting Trend

- Basic forecasting models for trends compensate for the lagging that would otherwise occur
- One model, **trend-adjusted exponential smoothing** uses a three step process
 - Step 1 - Smoothing the level of the series**
 - Step 2 - Smoothing the trend**
 - Forecast including the trend**

$$S_t = \alpha A_t + (1-\alpha)(S_{t-1} + T_{t-1})$$

$$T_t = \beta(S_t - S_{t-1}) + (1-\beta)T_{t-1}$$

$$FIT_{t+1} = S_t + T_t$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 670

Measuring Forecasting Accuracy

Mean Absolute Deviation (MAD)

➢ measures the total error in a forecast without regard to sign

$$MAD = \frac{\sum |actual - forecast|}{n}$$

Cumulative Forecast Error (CFE)

➢ Measures any bias in the forecast

$$CFE = \sum (actual - forecast)$$

Mean Square Error (MSE)

➢ Penalizes larger errors

$$MSE = \frac{\sum (actual - forecast)^2}{n}$$

Tracking Signal

➢ Measures if your model is working

$$TS = \frac{CFE}{MAD}$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 671

Stationarity

stationary time series have no trend.

conditions

1. constant mean
2. Constant variance
3. An autocovariance that does not depend on time

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 672

iid noise

The time series in which there is no trend or seasonal component and the observations are simply independent and identically distributed (iid) random variables with zero mean.

such sequence of random variables $x_1, x_2 \dots x_n$ as iid noise

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 673

Auto Correlation

$$\text{auto covariance}_h(x_t) = \text{cov}(x_t, x_{t-h})$$

$$\text{auto correlation}_h(x_t) = \frac{\text{Auto cov}_h(x_t)}{\text{std}(x_t) \text{std}(x_{t-h})}$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 674

Auto Correlation Function

$$\text{auto covariance}_h(x_t) \stackrel{\text{def}}{=} \text{cov}(x_t, x_{t-h})$$

$$\text{ACF} = \frac{\rho_x(h)}{\rho_x(0)} = \text{Cor}(x_t, x_{t-h})$$

$$\rho_x(h)$$

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 675

Models

- AR Model $\rightarrow AR(p)$
- MA Model $\rightarrow MA(q)$
- ARMA Model $\rightarrow ARMA(p,q)$
- ARIMA Model

Advanced Statistical Techniques for Analytics 26-04-2019 Slide 676

AR Model(Auto regressive model)

AR (p)

$$y_t = \delta_t + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \epsilon_t$$

Constant
value of time series
at time $t-j$

$$\phi_p \neq 0$$

$\epsilon_t \sim N(0, \sigma^2)$
for all t .

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 677

Moving Average(MA) Model

+ today's announcement
+ yesterday's

$$y_t = f(\epsilon_t, \epsilon_{t-1}, \epsilon_{t-2}, \dots)$$

$$MA(q) = \theta_0 + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

θ_k is constant for $k=1, 2, \dots, q$

$$\theta_0 \neq 0$$

$$\epsilon_t \sim N(0, \sigma^2) \text{ for all } t.$$

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 678

ARMA model –ARMA(p, q)

$$y_t = \delta + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

If $p \neq 0$ and $q=0$, then AR (p)

If $p=0$ and $q \neq 0$, then MA (q)

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 679

Selecting the Right Forecasting Model

1. The amount & type of available data
 - Some methods require more data than others
2. Degree of accuracy required
 - Increasing accuracy means more data
3. Length of forecast horizon
 - Different models for 3 month vs. 10 years
4. Presence of data patterns
 - Lagging will occur when a forecasting model meant for a level pattern is applied with a trend

Advanced Statistical Techniques for Analytics

26-04-2019

Slide 680