Knowledge Discovery & Data Mining Analyzing Feature Relationships Instructor: Yong Zhuang

yong.zhuang@gvsu.edu

Yong Zhuang

Outline

- Analyzing Feature Relationships
 - Introduction to Feature Analysis
 - Covariance (for numerical features)
 - Correlation Coefficient (for numerical features)
 - Spearman's Rank Correlation (Numeric & Ordinal Data)
 - Chi-Square Test (for categorical features)
 - Partial correlation



Why Analyze Relationships Between Features?

- Purpose: Understanding the relationships between features is key to improving predictive models, detecting patterns, and identifying significant associations.
- Key Reasons:
 - Identify Correlations: Determine how one feature may influence or be related to Ο another.
 - Improve Model Performance: Feature relationships can inform better feature Ο selection and model building.
 - Detect Patterns: Recognize trends and patterns within the data. Ο
 - Hypothesis Testing: Verify if observed relationships in the data are statistically Ο significant.





Covariance is measure assessing how much two attributes change together. Consider two numeric attributes A and B and a set of *n* real valued observations {(a1, b1), ..., (an, bn)}. The mean values (also known as the expected values) of A and B, that is,

$$E(A) = \bar{A} = \frac{\sum_{i=1}^{n} a_i}{n}$$

Then the covariance between A and B is defined as

$$Cov(A, B) = E((A - \overline{A})(B - \overline{B})) = \frac{\sum_{i=1}^{n} (a_i - \overline{A})(b_i - \overline{B})}{n}$$

or

 $Cov(A, B) = E(A \cdot B) - AB$

and
$$E(B) = \overline{B} = \frac{\sum_{i=1}^{n} b_i}{n}$$





- **Positive covariance**: If Cov(A,B) > 0, then A and B tend to increase together. As the values of A increase, the values of B also tend to increase, and similarly, as the values of A decrease, the values of B tend to decrease.
- **Negative covariance**: If Cov(A,B) < 0 then A and B tend to move in opposite directions. As the values of A increase, the values of B tend to decrease, and vice versa.

If A and B are independent, then $E(A \cdot B) = E(A) \cdot E(B)$. $\rightarrow Cov(A,B) = 0$







Example. This table presents a simplified example of stock prices observed at five time points for AllElectronics and HighTech, a high-tech company. If the stocks are affected by the same industry trends, will their prices rise or fall together?

- E(AllElectronics) =
- E(HighTech) =
- Cov(AllElectroncis, HighTech) =



Time point	AllElectronics	HighTech
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5

	Cov(A,	B) =	$E(A \cdot$	B)	$-\bar{A}$	B
--	--------	------	-------------	----	------------	---





the same industry trends, will their prices rise or fall together?

- E(AIIE ectronics) = 4
- E(HighTech) = 10.8
- Cov(AllElectroncis, HighTech)

 $= \frac{6 \times 20 + 5 \times 10 + 4 \times 14 + 3 \times 5 + 2 \times 5}{-4 \times 10.80}$

= 50.2 - 43.2 = 7.

Example. This table presents a simplified example of stock prices observed at five time points for AllElectronics and HighTech, a high-tech company. If the stocks are affected by

Time point	AllElectronics	HighTech
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5





Example. This table presents a simplified example of stock prices observed at five time points for AllElectronics and HighTech, a high-tech company. If the stocks are affected by the same industry trends, will their prices rise or fall together?

- E(AllElectronics) = 4
- E(HighTech) = 10.8
- Cov(AllElectroncis, HighTech)

$$=\frac{6 \times 20 + 5 \times 10 + 4 \times 14 + 3 \times 5 + 2 \times 5}{5} - 4 \times$$

= 50.2 - 43.2 = 7.

stock prices for both companies rise together

	Time point	AllElectronics	HighTech
	t1	6	20
	t2	5	10
	t3	4	14
10.80	t4	3	5
10.00	t5	2	5





Correlation Analysis (Numeric Data)

Correlation coefficient (also called Pearson's product moment coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^{n} (a_i - \bar{A})}{n\sigma_A \sigma}$$

- $(a_i b_i)$ is the sum of the AB cross-product.
- The higher the value, the stronger the correlation
- $r_{A,B} = 0$: independent; $r_{AB} < 0$: negatively correlated



where n is the number of tuples, \overline{A} and \overline{B} are the respective means of A and B, σ_A and σ_B are the respective standard deviation of A and B, and Σ

• $-1 \le r_{A,B} \le +1$, If $r_{A,B} > 0$, A and B are positively correlated (A's values increase as B's).



Visually Evaluating Correlation





0.40

0.60

0.50

0.70





1.00 0.90 0 DO 🕑

Scatter plots showing the correlation coefficient from –1 to 1.

Knowledge Discovery & Data Mining



Visually Evaluating Correlation





Knowledge Discovery & Data Mining

Correlation Coefficient (Numeric Data)

Example. $r_{AE,HT}$?

Cov(AE, HT) = 7



Yong Zhuang

 $r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$

Time point	AllElectronics	HighTech
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5







Correlation Coefficient (Numeric Data)

Example. $r_{AE,HT}$?

- Cov(AE, HT) = 7
- $\sigma_{AE} = \sqrt{2} \approx 1.414$
- $\sigma HT = \sqrt{32.56} \approx 5.706$

$r_{AE,HT} \approx 7/(1.414*5.706) \approx 0.868$

 $r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$

Time point	AllElectronics	HighTech
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5





$r_{AE,HT} \approx 0.868$

 $t = 0.87 \times \sqrt{\frac{5-2}{1-0.87^2}} = 0.87 \times \sqrt{\frac{3}{1-0.7569}} = 0.87 \times \sqrt{\frac{3}{0.2431}}$ $t = 0.87 \times \sqrt{12.35} = 0.87 \times 3.515 = 3.06$

features are correlated. **Correlated**



For 3 degrees of freedom, the t-statistic required to reject the null hypothesis at the 0.05 significance level is 2.353. Based on our computed t-statistic, we can reject the null hypothesis that AllElectronics and HighTech are independent, indicating that the two

Table of Critical Values for Student's t-Distribution.





Correlation Coefficient (Numeric Data)

Assumptions of Pearson Correlation

- All objects (data points) should be independent of each other.
- Ο relationships.
- Both attributes should be **continuous** and follow a **normal distribution**.

Limitations

• Pearson correlation is **extremely sensitive** to outliers.

The relationship between the two variables should be **linear**, not suitable for **nonlinear**







association between two ranked variables.

- monotonic function (variables move in the same direction).

Spearman's Rank Correlation Coefficient is a measure of the strength and direction of

• It is a non-parametric test, meaning it does not assume a normal distribution of the data. • It assesses how well the relationship between two variables can be described using a



Example.

Individual	Score in Test X	Score in Test Y
1	87	88
2	91	84
3	65	75
4	70	62
5	85	78

Knowledge Discovery & Data Mining



Example.

	Individual	Test X	Rank X	Test Y	Rank Y
1		87	4	88	5
2		91	5	84	4
3		65	1	75	2
4		70	2	62	1
5		85	3	78	3





Example.

	Individual	Rank X	Rank Y	d _i = Rank X - Rank Y	d _i ²
1	4		5	-1	1
2	5		4	1	1
3	1		2	-1	1
4	2		1	1	1
5	3		3	0	0

If no ties in the data

$$r_s = 1 - rac{6\sum d_i^2}{n(n^2-1)} = 1 - rac{6 imes(1+1+1+1)}{5 imes(5^2-1)} = 0$$

.8

Perfect positive monotonic relationship between Test A and Test B scores







Example with Tied Values.

Individual	Test X	Rank X	Test Y	Rank Y
1	82	3.5	88	5
2	82	3.5	84	4
3	65	1	75	2.5
4	70	2	62	1
5	85	5	75	2.5

 $r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$

Where A is Rank X and B is Rank Y





For nominal data, a correlation relationship between two attributes, A and B, can be discovered by a χ^2 (chi-square) test.



event has its own cell (or slot) in the table.

- Suppose A has c distinct values, namely a1,a2, ..., ac. B has r distinct values, namely
- b1,b2, ..., br. The data tuples described by A and B can be shown as a contingency
- table, with the c values of A making up the columns and the r values of B making up
- the rows. Let (Ai, Bj) denote the joint event that attribute A takes on value ai and attribute
- B takes on value bj, that is, where (A = ai, B = bj). Each and every possible (Ai, Bj) joint





The χ^2 value (also known as the Pearson χ^2 statistic) is computed as $\chi^{2} = \sum_{i=1}^{c} \sum_{j=1}^{r} \frac{(o_{ij} - e_{ij})^{2}}{e_{ij}}$ i = 1 j = 1

where oij is the observed frequency (i.e., actual count) of the joint event (Ai,Bj) and eij is

the expected frequency of (Ai,Bj), which can be computed as $e_{ij} = \frac{count(A)}{a}$

The cells that contribute the most to the χ^2 value are those whose actual count is very different from the expected count

$$= a_i) \times count(B = b_j)$$

n





he or she liked playing video games.



$$e_{ij}$$
 =

Yong Zhuang

- **Example.** Suppose that a group of 1500 people was surveyed. Each person was asked
- whether his or her preferred type of reading material was fiction or nonfiction, and whether

contingency table

No game(ng)	Sum (row)
200	
1000	



Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether he or she liked playing video games.

contingency table

	Game(g)	No game(ng)	Sum (row)
Fiction(f)	250	200	?
Nonfiction(nf)	50	1000	?
Sum(col.)	?	?	?



Knowledge Discovery & Data Mining



Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether he or she liked playing video games.

contingency table

	Game(g)	No game(ng)	Sum (row)
Fiction(f)	250	200	450
Nonfiction(nf)	50	1000	1050
Sum(col.)	300	1200	1500



Knowledge Discovery & Data Mining



Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether he or she liked playing video games.

contingency table

	Game(g)	No game(ng)	Sum (row)
Fiction(f)	250 (e f,g ?)	200(ef,ng?)	450
Nonfiction(nf)	50(enf,g?)	1000(e nf,ng?)	1050
Sum(col.)	300	1200	1500
$e_{ij} = -$	count(A =	$= a_i) \times count($	$(B = b_j)$

n

Knowledge Discovery & Data Mining



Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether he or she liked playing video games.

contingency table

	Game(g)	No game(ng)	Sum (row)
Fiction(f)	250 (90)	200(e f,ng?)	450
Nonfiction(nf)	50(e nf,g ?)	1000(enf,ng?)	1050
Sum(col.)	300	1200	1500

 $e_{f,g} = \frac{300 \times 450}{1500} = 90$





Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether he or she liked playing video games.

contingency table

	Game(g)	No game(ng)	Sum (row)
Fiction(f)	250 (90)	200(360)	450
Nonfiction(nf)	50(210)	1000(840)	1050
Sum(col.)	300	1200	1500



Knowledge Discovery & Data Mining



he or she liked playing video games.

contingency table



Example. Suppose that a group of 1500 people was surveyed. Each person was asked whether his or her preferred type of reading material was fiction or nonfiction, and whether

$$\chi^{2} = \sum_{i=1}^{c} \sum_{j=1}^{r} \frac{(o_{ij} - e_{ij})^{2}}{e_{ij}}$$

No game(ng)	Sum (row)
200(360)	450
1000(840)	1050
1200	1500

$$\frac{200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840}$$
$$8 = 507.93.$$





$\chi^2 = 507.93$

the χ^2 distribution).

Based on our computed value, we can reject the hypothesis that play game and preferred_reading are independent, so they are **Correlated**.

The χ^2 statistic tests the hypothesis that A and B are independent, that is, there is no correlation between them. The test is based on a significance level, with $(r - 1) \times (c - 1)$ degrees of freedom. Since in this example, r = 2 and c = 2, the degrees of freedom are (2) -1×(2 – 1) = 1. For 1 degree of freedom, the χ^2 value needed to reject the hypothesis at the 0.001 significance level is 10.828 (taken from the table of upper percentage points of

Table of upper percentage points of the Chi-squared distribution

Knowledge Discovery & Data Mining

Partial correlation

Suppose we have the following data on three variables, X, Y, and Z:

X	Y	z
2	1	0
6	3	0
10	2	1
20	4	1

These data have the feature that whenever Z = 0, X equals exactly twice Y, and whenever Z = 1, X is exactly 5 times Y. Thus, contingent on the value of Z, there is an exact relationship between X and Y; but the relationship cannot be said to be exact without reference to the value of Z.

In fact, if we compute the Pearson correlation coefficient between variables *X* and *Y*, the result is approximately 0.836, while if we compute the partial correlation between *X* and *Y*, using the formula given below, we find a partial correlation of 0.919, which is stronger than the full correlation.





Partial correlation

$$egin{aligned} \mathbf{w}_X^* &= rg\min_{\mathbf{w}} \left\{ \sum_{i=1}^N (x_i - \langle \mathbf{w}, \mathbf{z}_i
angle)^2
ight\} \ \mathbf{w}_Y^* &= rg\min_{\mathbf{w}} \left\{ \sum_{i=1}^N (y_i - \langle \mathbf{w}, \mathbf{z}_i
angle)^2
ight\} \end{aligned}$$

with N being the number of observations and $\langle \mathbf{w}, \mathbf{v} \rangle$ the scalar product between the vectors w and v.

The residuals are then

$$egin{aligned} e_{X,i} &= x_i - \langle \mathbf{w}_X^*, \mathbf{z}_i
angle \ e_{Y,i} &= y_i - \langle \mathbf{w}_Y^*, \mathbf{z}_i
angle \end{aligned}$$

and the sample partial correlation is then given by the usual formula for sample correlation, but between these new derived values:

$$\hat{\rho}_{XY\cdot\mathbf{Z}} = \frac{N\sum_{i=1}^{N} e_{X,i} e_{Y,i} - \sum_{i=1}^{N} e_{X,i} \sum_{i=1}^{N} e_{Y,i}}{\sqrt{N\sum_{i=1}^{N} e_{X,i}^2 - \left(\sum_{i=1}^{N} e_{X,i}\right)^2} \sqrt{N\sum_{i=1}^{N} e_{Y,i}^2 - \left(\sum_{i=1}^{N} e_{Y,i}\right)^2}}.$$

A simple way to compute the sample partial correlation for some data is to solve the two associated linear regression problems, get the residuals, and calculate the correlation between the residuals. Let X and Y be, as above, random variables taking real values, and let Z be the *n*-dimensional vector-valued random variable. We write x_i , y_i and z_i to denote the *i*th of N i.i.d. observations from some joint probability distribution over real random variables X, Y and Z, with z_i having been augmented with a 1 to allow for a constant term in the regression. Solving the linear regression problem amounts to finding (n+1)-dimensional regression coefficient vectors \mathbf{w}_X^* and \mathbf{w}_Y^* such that

Knowledge Discovery & Data Mining



Summary

- Feature Analysis: Relationships
 - Introduction to Feature Analysis
 - Covariance (for numerical features)
 - Correlation Coefficient (for numerical features)
 - Spearman's Rank Correlation (Numeric & Ordinal Data)
 - Chi-Square Test (for categorical features)
 - Partial correlation

