## دانشگاه ملی مهارت

## آموزشكده ميناب

نام و نام خانوادگی: آرش زارعیان

واحد درسى: مباحث ويژه

رشته: مهندسی حرفه ای کامپیوتر

مدرس: محمد احمد زاده

## : Data Preprocessing 4 بخش

- A. چرا Data Cleaning در علم داده اهمیت دارد؟
  - Missing Values .B چگونه مدیریت می شوند؟
- Outliers .Cچیست و چگونه می توانید آنها را تشخیص دهید؟
  - Data Transformation .D چرا کاربرد دارد؟
- Encoding Techniques (One-Hot Encoding) جه تفاوتي دارند؟
  - Feature Selection در Model-building اهمیت دارد؟ .
    - Duplicate Data .Gچگونه در پایگاه دادهها حذف می شود؟
  - Irrelevant Data .Hچه مشکلاتی را در پیشبینیهای Machine Learning ایجاد می کند؟
    - ا. چرا Data Imputation برای پر کردن Missing Values کاربرد دارد؟
      - J. چگونه می توانید Normality را در دادههای عددی بررسی کنید؟

	6
Day. Month. Year. Subject. Sys - 16	
DayMonthYear Subject	7
Data Proprocessing 4 curs	2
Posta deaning 17: A	3
Data cleaning بلى ازم توبن مراحل دربودازش وتعلى طده ها است . دلایل	5
9,000,000	6
العيب كن عيارت الداز:	-7
ارافرانس وقت سرل ها ٢ ماهش فورونمافهات ٣ بهود عمارد اللورسم ها	9
ا معدم لیری میں کہ را اور کی در زمان ر هریند .	
م طور ملی ، لیفیت داده ها ما نمو صدیقتی بر سایج مطلبل و سرک های بادلوی مانسی دارد.	13
	7.4
بعرق بالمعازى مناسب، صفى فوى ترين اللورية ها ينزجى و ابنز عملار مطلوبي طالبية بالله	15
Softs Tune Wissing Values B	16
MINES 12	17
مدرس مقادیر کے بشری در وادہ ہا بخس مہی از بسو پردازیں دارہ ہا است روس ہای	18
Control of the Contro	19
معالی برای براور را این معادی و اور دارد ام مستم بدنوع داره و کاربرد.	21
	_22_
المتقاب مي ويون	-23
	_24

DayMonthYear	Subject.
	ررتو) های موبویت معادیر کے بشرہ
ولار مع على في عمل العب	مرسون داده های دارای منوار کم نسره مزیت: ساده
	عت (زدست رقبق) الحلاعات مهر لنرج و
بیب: صالی است دفت مدل ماهنس ما م	المرتبي عرب : دادهها لفظ ي لوند ع
د فلو کا عیب : زمان بردازتی با	۲ اسفاده از درل های بسسنی مربت: مرابلزین
	ع السفاده ازمغدار قبلی ما بعدی مرست: مراس
ست داده ها مادرست سویز	عيب: در هورت و او حقيرات فالهادي ممل اه
فريت : مارهسي نيازيد يوسي بردارس	ک. اسفاده از اللورین های مقادم بر مقادیر کی سره
	Time just jund ; wie
	ا مراز المارات و معلونه ي حوال المارار
	out liers به نعاط داده ای دست می انود کرم طور قابل در
سَلَفَى ما يَسْرُكُوا هاى إندازه ليرك	صوعه داده دفاوت دارند ابن نقاط مي آيند به دلايل م
	AVANGE

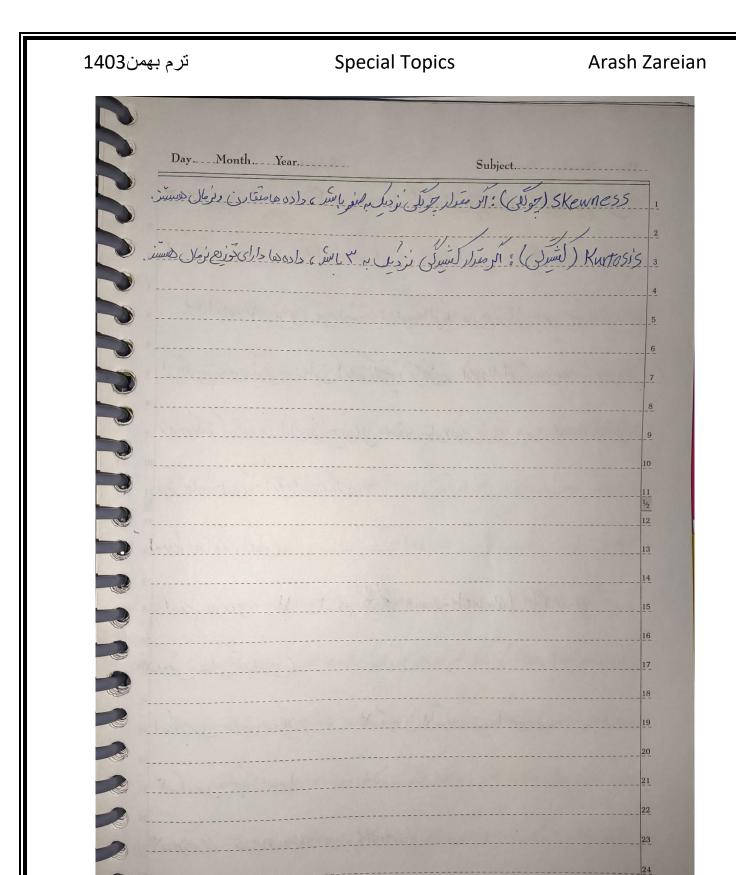
DayMonthYear	Subject.
	Subject.
ه ها در ور ساس راه های صلعی سرای	ورود انستان جاده ها یا بربره دمای واقعی در داد
	ز. بساسان داده های برت و اور دارد از علم ه
۲ استاده از انوان مار	اران اماری (مصدرده من میاری - IOR )
٢ / ١١١١ عاده / ر تيو د ار دی	کروس های مسنی بر معل های مادلیری طالسی
All and the state of the state	Pata transformation > \$
12 (20) 4 Jan - 1811 1916	راس سای Data Transformation
chi il ser chi (it). Sollela osts	است درای بهبود معلی، سل سازی رفسیر
، رسم سری ، رم لذاری و ماهسی العاد ما سر	( سازار دسازی ، نرمال سازی ، سیل ساریی
ع الدرى مالس ٢ مدريت داده هاى رب	اربردهای مهم: ا-لبود دنت ملهای
مربهور توزيع داده ها كى كاهش ابعا	٧ ـ ا مزاس كاراني اللورسة هاى مسنى بر فاصله
ر معلف معتلف معتلف	ميل داده ها م سازماری داده ها ساز
Lable Lack to Lable	Encoding, one-Hot Encoding E

9	
)	D W. J. W.
	DayMonthYear
	ا فرد روس برای سریل متعلم های داسترسی شره بر داده های عددی (سیفاده می امون ایا
	1 - 2 < 6 (616) 1615 MICH (2015) - MAINS MICHAEL A SCENE
	و خوه کاربرد و تسویل اکرنامیفارت (ست.
	Canol Line Sill Find Singer 1
	(Could (Silver) table Encoding 5
8	م دراین دوس هر داست ( للاس) به بکر عدد محمد احتصاص داده ج التوج
) .	و مرایا : ساخی و کاهش ابعاد داده ها - صاسب برای سل هایی له روابط تر سی وادر نقل می کود
-	و مرایا : سادی و ناهس ابعاد داده ها - صافست برای سدل هایی که روابط بر نبعی وادر طرح مراید
3	10
	المعاس: الحاد ارتباط هاى) عدد نادرس سن کسترها
3	ال بعایب: ایجاد ارتباط های عدد نادرست سن دلسته ها الم
•	(30 Li C19) UHE-401 PHOCHY 13
	The state of the s
	15 حراین روسی برای هر مقدار از صفیر داستر بنری سری میر ویژگی محدار ادراد سری و مقدار ۱ یا و
3	16
) -	١٦ مراي ( ونعالي ي يايد
	20 C C ( C ( C ) ( ) 17 17
3	18
رُن ا	١٥ سرايا: ازسي بردن ارساط عددي فاحرست بين ولسيرها - مناسب براي مول هادي موض
2	
-	ورده های ورودی مسیقل از بلدبلرین
3) 1-	21 ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) ( ) (
2 -	22
	23 معاید: افزانس ابعاد داده ، به کیوس در صفیرهای در سری سری سری ما معادی زیار
	24
-	
=	
2	

i D. M. d. V.
Day. Month. Year. Subject.
Sistem 1 Model-building is Feature selection 15 F
Feature Selection فراسر استاب مهرترین در این ا زمیان بدای ویز لی های موجود
5 - در صفر عد داده (بست . این اور ماعت به دو کادایی سال عما دیس سعیدی ر جار لیری (ز
i) Interesting over litting 7
e I ) San July Ovar fitting of the early of 1
الله مر کلفتن زع ن مردازش و هزید محاسات کر بهبود دفیسر پذیری صل
الم من بالما على ماده ها المزون المود. ور بالما على ماده ها المزون المود. الما على المده ها المزون المود.
داده های مالای زبانی درجی دهند در رون های مسام یا باسای دربانها و داده در در ه سوند.
از معرف این داده ها برای بمبود کارای ، کاهس مع داده ها و کلولوی از بحلیل های نادرس موری ال
اه اور روستی ها: ار استاده از SQL برای کرد دادههای مارای
Pandas 2 (c) Ju (c) 21 21
Exect > 6 (1) by (1) (22 )
23
EVANGE -

	DayMonthYear	Subject	
9 july	Machine Learning (Suscing	الريسيلاتي رادرسية Irrelevant Data . ا	11
		الرداده ها ى نامرجوط) مرويزكي م المريزكي م	2
	ست ماعت كالملكي دفت مدل كنوين	ى مدل بادلىرى مائسن ارائدىنى ھەند و محمى صلى ا	5 برا
	/ کارایس) افز اس پیپیرلی و ایجاد	بن نوع حاده ها معرلاً نون اصاد لرده و باعث المصر	7
	ر السال رو في السال .	مرادر بر مسلات به وای را	7 9
(-		ا کا کاسی دفت مدک ، ۲ افزانس وفت	12
)	ر سل های مینی بر فاهدار	المسترسين تفسير صل که الاهستي عملاً	13
	Missing Value	es issur Sly Pata Imputation Jr. I	15
رية المراق		Data ImPutation بو مانیک مے دریس بردارند	
	ی از داده های از داست راسم مانسر منع	سقولهای صاسب ما محایلترین می کند . این کار ماعث	19
	ل ودن معري داده لمال نظر ,	روی دفت مدل نعانسته مانشد و مرافقط لیست و ماما	_21
ر مساری	مع داده ها کرد بهبود دفت مدل بارلیم	دلایل اهمیت در ارساد : استلولس از مادسی	_22_
)			24

DayMonthYear	Subject
ع افزایس بایراری رقابلیت تعمیم مدل	1 سے الولای از محافظ مر داده ها
	Mormality july Jest 3
ه ها از بل درنع نرمال سردی می لیند ساری	5 Normali47/ 5
( ماری المرسل A NOVA رکون های تطاری بارامتریل )	۲ - از اللوريع های بادليری مالسي (مانيز دلرسور)
مر مربرای بررسی برمال جودن داده ها می قران از	و فرور) می استر کرده ها دارای توزیع سرمال هسی
	اه در در ایرای (سعفا ده تورد) ا در در در ایرای (سعفا ده تورد) ا در
	12 13 - 10 m ( 12 ) h ( 12 )
Q-Q Plot (Quantile-Quanti	le (1) -Y Histogram (1)-1 15
	16 (Sla) (Sla) 17
ن لولمولروف _ اسريروف	اهر ار را مون سا سرد - وبلا کرانور انور انور انور انور انور انور انور
	20 کے کر زمون سارک مدیرا 21 کے کر زمون سارک مدیرا
Ke	irtosis , skewness (syr-723
	24



EVANGE