Pandas 資料分析 (2)

張家瑋 副教授 國立臺中科技大學資訊工程系

選取多個 DataFrame 的欄位

In [10]:		ort pandas a ort numpy as										
			d_csv('data/mo ector = movies	ovie.csv') s[['actor_1_name'	, 'actor_2_name	e',						
	mov:	ie_actor_dir	ector.head()	'actor_3_name'	, 'director_nam	ne']]						
Out[10]:		actor_1_name	actor_2_name	actor_3_name	director_name							
	0		Joel David Moore	Wes Studi								
	1		Orlando Bloom		Gore Verbinski							
	2	Christoph Waltz										
	3	Tom Hardy	Christian Bale	Joseph Gordon-Levitt	Christopher Nolan							
	4	Doug Walker	Rob Walker	NaN	Doug Walker							
To [11].	+1100	مرسمین محالا ط	irector name'	11\								
		,		11)								
Out[11]:	pan	das.core.fra	me.DataFrame									
In [12]:	typ	e(movies['di	rector_name'])								
Out[12]:	pan	das.core.ser	ies.Series									
- [40]				1225								
		•	[:, ['director	_name[]])								
Out[13]:	pan	das.core.fra	me.DataFrame									
In [14]:	typ	e(movies.loc	[:, 'director	_name'])								
Out[14]:	pan	das.core.ser	ies.Series									
In [15]:	5]: cols = ['actor_1_name', 'actor_2_name', 'actor_3_name', 'director_name'] movie_actor_director = movies[cols]											

	actor_1_name	actor_2_name	actor_3_name	director_name
0	CCH Pounder	Joel David Moore	Wes Studi	James Cameron
1	Johnny Depp	Orlando Bloom	Jack Davenport	Gore Verbinski
2	Christoph Waltz	Rory Kinnear	Stephanie Sigman	Sam Mendes
3	Tom Hardy	Christian Bale	Joseph Gordon-Levitt	Christopher Nolan
4	Doug Walker	Rob Walker	NaN	Doug Walker
4911	Eric Mabius	Daphne Zuniga	Crystal Lowe	Scott Smith
4912	Natalie Zea	Valorie Curry	Sam Underwood	NaN
4913	Eva Boehnke	Maxwell Moody	David Chandler	Benjamin Roberds
4914	Alan Ruck	Daniel Henney	Eliza Coupe	Daniel Hsia
4915	John August	Brian Herzlinger	Jon Gunn	Jon Gunn

用方法 (Methods) 選取欄位 (1/3)

```
movies = pd.read csv('data/movie.csv')
movies.dtypes.value counts()
float64
              13
object
              12
 int64
               3
dtype: int64
movies.select dtypes(include='object').head()
    color director name actor 2 name
                                                              genres actor 1 name movie title actor 3 name
                                                                                                                                         plot keywords
                   James
                               Joel David
                                          Action|Adventure|Fantasy|Sci-
                                                                       CCH Pounder
                                                                                                       Wes Studi
                                                                                                                      avatar|future|marine|native|paraplegic
 0 Color
                                                                                          Avatar
                                   Moore
                 Cameron
                                                                                        Pirates of
                                                                                                           Jack
                                                                                                                      goddess|marriage ceremony|marriage http://www.ii
           Gore Verbinski Orlando Bloom
                                              Action|Adventure|Fantasy
                                                                        Johnny Depp
                                                                                      Caribbean:
                                                                                                      Davenport
                                                                                                                                           proposal|pi.
                                                                                       At World's
                                                                                            End
                                                                           Christoph
                                                                                                       Stephanie
                                                                                                                                                         http://www.ii
                            Rory Kinnear
                                               Action|Adventure|Thriller
                                                                                         Spectre
                                                                                                                       bomb|espionage|sequel|spy|terrorist
 2 Color
             Sam Mendes
                                                                                                        Sigman
                                                                               Waltz
                                                                                        The Dark
               Christopher
                                                                                                                 deception|imprisonment|lawlessness|police http://www.ii
 3 Color
                            Christian Bale
                                                         Action|Thriller
                                                                          Tom Hardy
                                                                                          Knight
                                                                                                   Gordon-Levitt
                    Nolan
                                                                                           Rises
                                                                                       Star Wars:
                                                                                                                                                        http://www.ii
                                                                                      Episode VII
     NaN
             Doug Walker
                              Rob Walker
                                                         Documentary
                                                                        Doug Walker
                                                                                                           NaN
                                                                                      - The Force
                                                                                        Awakens
```

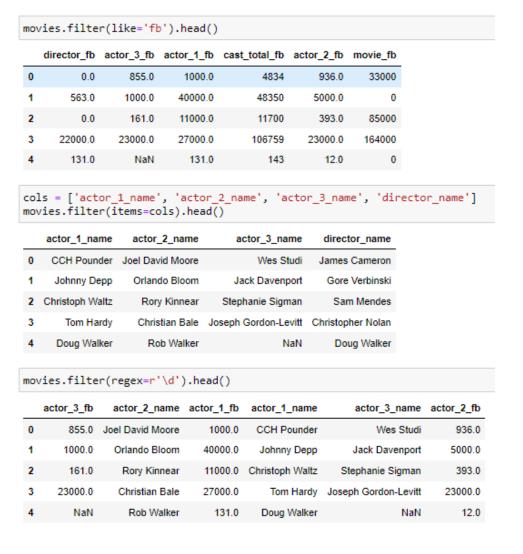
用方法 (Methods) 選取欄位 (2/3)

	num_c	ritic dura	ition	director_fb	actor_3_fb	actor_1_fb	gro	oss num_vot	ed_users c	ast_total_fb	facenumber_	in_poster	num_user	budg	et title_y
0	7:	23.0 1	78.0	0.0	855.0	1000.0	76050584	7.0	886204	4834		0.0	3054.0	237000000	.0 200
1	3	02.0 1	69.0	563.0	1000.0	40000.0	30940415	2.0	471220	48350		0.0	1238.0	300000000	.0 200
2	6	02.0 1	48.0	0.0	161.0	11000.0	20007417	5.0	275868	11700		1.0	994.0	245000000	.0 201
3	8	13.0 1	64.0	22000.0	23000.0	27000.0	44813064	2.0	1144337	106759		0.0	2701.0	250000000	.0 201
4	ı	NaN	NaN	131.0	NaN	131.0	N	laN	8	143		0.0	NaN	Na	N N
															-
10\			•	(include= actor_2_na		object']).		actor_1_name	movie_title	e actor_3_na	me		plot_	keywords	
0		director_r	name	actor_2_nar	me vid Action A	object']).	genres a	actor_1_name	movie_title			ıtar future m			nttp://www.i
	color	director_r	name ames neron	actor_2_nar	me vid Action <i>I</i> ore		genres a lasy Sci- Fi			f Wes Si	tudi ava ack gode		arine native	paraplegic	•
0	color	director_r Ja Can	name ames neron binski	actor_2_nai	me vid Action A ore om Act	dventure Fant	genres a asy Sci- Fi Fantasy	CCH Pounder Johnny Depp Christoph	Avatar Pirates of the Caribbean:	f J Daven	ack gode	dess marria	arine native	paraplegic / marriage /posal pi	•
0	Color Color	director_r J: Can Gore Vert	name ames neron binski	Joel Da Mod	me vid Action / ore Act om Act	dventure Fant ion Adventure tion Adventure	genres a asy Sci- Fi Fantasy	CCH Pounder Johnny Depp	Avatar Pirates of the Caribbean At World's End	f J Daven; Stepha Sign	ack gode port gode port bo	dess marria omb espiona	ge ceremony	paraplegic //marriage pposal/pi	http://www.i

用方法 (Methods) 選取欄位 (3/3)

•										
noι	/ies.filte	r(like='f	o').head())						
	director_fb	actor_3_fb	actor_1_fb	cast_total_fb	actor_2_fb	movie_fb				
0	0.0	855.0	1000.0	4834	936.0	33000				
1	563.0	1000.0	40000.0	48350	5000.0	0				
2	0.0	161.0	11000.0	11700	393.0	85000				
3	22000.0	23000.0	27000.0	106759	23000.0	164000				
4	131.0	NaN	131.0	143	12.0	0				
<pre>cols = ['actor_1_name', 'actor_2_name', 'actor_3_name', 'director_name'] movies.filter(items=cols).head()</pre>										
0	actor_1_na	der Joel Dav	_2_name	actor_3_na Wes S		tor_name Cameron				
1	Johnny De		do Bloom	Jack Daven		Verbinski				
2	Christoph Wa		y Kinnear	Stephanie Sign		m Mendes				
3	Tom Ha			oseph Gordon-Le		her Nolan				
4	Doug Wal	•	b Walker	•		ug Walker				
mo\	/ies.filte	r(regex=r	'\d').head	i()						
	actor_3_fb	actor_2_n	ame actor_	_1_fb actor_1	_name	actor_3_na	me actor_2			
0	855.0	Joel David M	loore 10	000.0 CCH P	ounder	Wes St	tudi 93			
1	1000.0	Orlando B	loom 400	000.0 Johnn	у Dерр	Jack Davenp	port 500			
2	161.0	Rory Kin	near 11	000.0 Christopl	n Waltz S	tephanie Sign	nan 39			
3	23000.0	Christian	Bale 27	000.0 Tom	Hardy Jose	ph Gordon-Le	evitt 2300			

用方法 (Methods) 選取欄位 (3/3)



在三種參數是互斥關係, 次只能用

在select_dtypes使用exclude排除欄位

mo	vies.s	elect_dtypes	(exclude='fl	oat').head()						
	color	director_name	actor_2_name	genres	actor_1_name	movie_title	num_voted_users	cast_total_fb	actor_3_name	
0	Color	James Cameron	Joel David Moore	Action Adventure Fantasy Sci- Fi	CCH Pounder	Avatar	886204	4834	Wes Studi	avatar future r
1	Color	Gore Verbinski	Orlando Bloom	Action Adventure Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	48350	Jack Davenport	goddess marri
2	Color	Sam Mendes	Rory Kinnear	Action Adventure Thriller	Christoph Waltz	Spectre	275868	11700	Stephanie Sigman	bomblespior
3	Color	Christopher Nolan	Christian Bale	Action Thriller	Tom Hardy	The Dark Knight Rises	1144337	106759	Joseph Gordon-Levitt	deception imprisoni
4	NaN	NaN Doug Walker Rob Walker Document		Documentary	Doug Walker	Star Wars: Episode VII - The Force Awakens	8	143	NaN	
4										+

在select_dtypes使用exclude排除欄位

mo	vies.s	elect_dtypes	(exclude='fl	oat').head()						
	color	director_name	actor_2_name	genres	actor_1_name	movie_title	num_voted_users	cast_total_fb	actor_3_name	
0	Color	James Cameron	Joel David Moore	Action Adventure Fantasy Sci- Fi	CCH Pounder	Avatar	886204	4834	Wes Studi	avatar future
1	Color	Gore Verbinski	Orlando Bloom	Action Adventure Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	471220	48350	Jack Davenport	goddess marr
2	Color	Sam Mendes	Rory Kinnear	Action Adventure Thriller	Christoph Waltz	Spectre	275868	11700	Stephanie Sigman	bomb espio
3	Color	Christopher Nolan	Christian Bale	Action Thriller	Tom Hardy	The Dark Knight Rises	1144337	106759	Joseph Gordon-Levitt	deception imprison
4	NaN	Doug Walker	Rob Walker	Documentary	Doug Walker	Star Wars: Episode VII - The Force Awakens	8	143	NaN	
4										→

對欄位名稱進行排序 (1/2)

對欄位名稱進行排序 (2/2)

True

movies[new_col_order].head()

	movie_title	title_year	content_rating	genres	director_name	actor_1_name	actor_2_name	actor_3_name	color	country	 movie_fb
(0 Avatar	2009.0	PG-13	Action Adventure Fantasy Sci- Fi	James Cameron	CCH Pounder	Joel David Moore	Wes Studi	Color	USA	 33000
	Pirates of the 1 Caribbean: At World's End	2007.0	PG-13	Action Adventure Fantasy	Gore Verbinski	Johnny Depp	Orlando Bloom	Jack Davenport	Color	USA	 0
:	2 Spectre	2015.0	PG-13	Action Adventure Thriller	Sam Mendes	Christoph Waltz	Rory Kinnear	Stephanie Sigman	Color	UK	 85000
;	The Dark 3 Knight Rises	2012.0	PG-13	Action Thriller	Christopher Nolan	Tom Hardy	Christian Bale	Joseph Gordon-Levitt	Color	USA	 164000
•	Star Wars: Episode VII - The Force Awakens	NaN	NaN	Documentary	Doug Walker	Doug Walker	Rob Walker	NaN	NaN	NaN	 0

5 rows × 28 columns

DataFrame的統計方法 (1/3)

movies = pd.read_csv('data/movie.csv')
movies.shape

(4916, 28)

movies.size 137648 列數和行數的乘積

ndim=2 為 Dataframe ndim=1 為 Series

len(movies)
4916 列數為資料筆數

color 4897 4814 director name num critic for reviews 4867 duration 4901 director facebook likes 4814 actor 3 facebook likes 4893 actor 2 name 4903 actor 1 facebook likes 4909 gross 4054 4916 genres actor 1 name 4909 4916 movie title num voted users 4916 cast total facebook likes 4916 actor 3 name 4893 facenumber in poster 4903 plot keywords 4764 movie imdb link 4916 num user for reviews 4895 4904 language 4911 country content rating 4616 budget 4432 4810 title year actor_2_facebook_likes 4903 imdb score 4916 4590 aspect ratio movie_facebook_likes 4916 dtype: int64

movies.count()

movies.min() num critic for reviews duration director facebook likes actor 3 facebook likes actor_1_facebook_likes 162 gross Action genres movie title #Horror num voted users cast_total_facebook_likes facenumber in poster http://www.imdb.com/title/tt0006864/?ref =fn t... movie imdb link num user for reviews budget 218 title year 1916 actor 2 facebook likes imdb score 1.6 1.18 aspect_ratio movie_facebook_likes dtype: object

其他統計方法還有 max(), mean(), median(), std()

DataFrame的統計方法 (2/3)

ovies.describe().T								
	count	mean	std	min	25%	50%	75%	max
num_critic_for_reviews	4867.0	1.379889e+02	1.202394e+02	1.00	49.00	108.00	191.00	8.130000e+02
duration	4901.0	1.070908e+02	2.528602e+01	7.00	93.00	103.00	118.00	5.110000e+02
director_facebook_likes	4814.0	6.910145e+02	2.832954e+03	0.00	7.00	48.00	189.75	2.300000e+04
actor_3_facebook_likes	4893.0	6.312763e+02	1.625875e+03	0.00	132.00	366.00	633.00	2.300000e+04
actor_1_facebook_likes	4909.0	6.494488e+03	1.510699e+04	0.00	607.00	982.00	11000.00	6.400000e+05
gross	4054.0	4.764451e+07	6.737255e+07	162.00	5019656.25	25043962.00	61108412.75	7.605058e+08
num_voted_users	4916.0	8.264492e+04	1.383222e+05	5.00	8361.75	33132.50	93772.75	1.689764e+06
cast_total_facebook_likes	4916.0	9.579816e+03	1.816432e+04	0.00	1394.75	3049.00	13616.75	6.567300e+05
facenumber_in_poster	4903.0	1.377320e+00	2.023826e+00	0.00	0.00	1.00	2.00	4.300000e+01
num_user_for_reviews	4895.0	2.676688e+02	3.729348e+02	1.00	64.00	153.00	320.50	5.060000e+03
budget	4432.0	3.654749e+07	1.002427e+08	218.00	6000000.00	19850000.00	43000000.00	4.200000e+09
title_year	4810.0	2.002448e+03	1.245398e+01	1916.00	1999.00	2005.00	2011.00	2.016000e+03
actor_2_facebook_likes	4903.0	1.621924e+03	4.011300e+03	0.00	277.00	593.00	912.00	1.370000e+05
imdb_score	4916.0	6.437429e+00	1.127802e+00	1.60	5.80	6.60	7.20	9.500000e+00
aspect_ratio	4590.0	2.222349e+00	1.402940e+00	1.18	1.85	2.35	2.35	1.600000e+01
movie_facebook_likes	4916.0	7.348294e+03	1.920602e+04	0.00	0.00	159.00	2000.00	3.490000e+05

DataFrame的統計方法 (3/3)

movies.describe(percentiles=[.99]).T

	count	mean	std	min	50%	99%	max
num_critic_for_reviews	4867.0	1.379889e+02	1.202394e+02	1.00	108.00	5.466800e+02	3.130000e+02
duration	4901.0	1.070908e+02	2.528602e+01	7.00	103.00	1.890000e+02	5.110000e+02
director_facebook_likes	4814.0	6.910145e+02	2.832954e+03	0.00	48.00	1.600000e+04	2.300000e+04
actor_3_facebook_likes	4893.0	6.312763e+02	1.625875e+03	0.00	366.00	1.100000e+04	2.300000e+04
actor_1_facebook_likes	4909.0	6.494488e+03	1.510699e+04	0.00	982.00	4.492000e+04	6.400000e+05
gross	4054.0	4.764451e+07	6.737255e+07	162.00	25043962.00	3.264128e+08	7.605058e+08
num_voted_users	4916.0	8.264492e+04	1.383222e+05	5.00	33132.50	6.815846e+05	1.689764e+06
$cast_total_facebook_likes$	4916.0	9.579816e+03	1.816432e+04	0.00	3049.00	6.241390e+04	6.567300e+05
facenumber_in_poster	4903.0	1.377320e+00	2.023826e+00	0.00	1.00	8.000000e+00	4.300000e+01
num_user_for_reviews	4895.0	2.676688e+02	3.729348e+02	1.00	153.00	1.999240e+03	5.060000e+03
budget	4432.0	3.654749e+07	1.002427e+08	218.00	19850000.00	2.000000e+08	4.200000e+09
title_year	4810.0	2.002448e+03	1.245398e+01	1916.00	2005.00	2.016000e+03	2.016000e+03
actor_2_facebook_likes	4903.0	1.621924e+03	4.011300e+03	0.00	593.00	1.700000e+04	1.370000e+05
imdb_score	4916.0	6.437429e+00	1.127802e+00	1.60	6.60	8.500000e+00	9.500000e+00
aspect_ratio	4590.0	2.222349e+00	1.402940e+00	1.18	2.35	4.000000e+00	1.600000e+01
movie_facebook_likes	4916.0	7.348294e+03	1.920602e+04	0.00	159.00	9.385000e+04	3.490000e+05

DataFrame的統計方法 – skipna = False

movies.min()	
num_critic_for_reviews	1
duration	7
director_facebook_likes	0
actor_3_facebook_likes	0
actor_1_facebook_likes	0
gross	162
genres	Action
movie_title	#Horror
num_voted_users	5
cast_total_facebook_likes	0
facenumber_in_poster	0
movie_imdb_link	http://www.imdb.com/title/tt0006864/?ref_=fn_t
num_user_for_reviews	1
budget	218
title_year	1916
actor_2_facebook_likes	0
imdb_score	1.6
aspect_ratio	1.18
movie_facebook_likes	0
dtype: object	

<pre>movies.min(skipna=False)</pre>	
num_critic_for_reviews	NaN
duration	NaN
director_facebook_likes	NaN
actor_3_facebook_likes	NaN
actor_1_facebook_likes	NaN
gross	NaN
genres	Action
movie_title	#Horror
num_voted_users	5
cast_total_facebook_likes	0
facenumber_in_poster	NaN
movie_imdb_link	http://www.imdb.com/title/tt0006864/?ref_=fn_t
num_user_for_reviews	NaN
budget	NaN
title_year	NaN
actor_2_facebook_likes	NaN
imdb_score	1.6
aspect_ratio	NaN
movie_facebook_likes	0
dtype: object	

С	olor	director_name	num_critic	duration	director_fb	actor_3_fb	actor_2_name	actor_1_fb	gross	genres	 num_user	language	country	content_ratin
0 F	alse	False	False	False	False	False	False	False	False	False	 False	False	False	Fals
1 F	alse	False	False	False	False	False	False	False	False	False	 False	False	False	Fals
2 F	alse	False	False	False	False	False	False	False	False	False	 False	False	False	Fals
3 F	alse	False	False	False	False	False	False	False	False	False	 False	False	False	Fals
4	True	False	True	True	False	True	False	False	True	False	 True	True	True	Tru

5 rows × 28 columns

15

movies.isna().sum().sum()

2654

color	True
director_name	True
num_critic	True
duration	True
director_fb	True
actor_3_fb	True
actor_2_name	True
actor_1_fb	True
gross	True
genres	False
actor_1_name	True
movie_title	False
num_voted_users	False
cast_total_fb	False
actor_3_name	True
facenumber_in_poster	True
plot_keywords	True
movie_imdb_link	False
num_user	True
language	True
country	True
content_rating	True
budget	True
title_year	True
actor_2_fb	True
imdb_score	False
aspect_ratio	True
movie_fb	False
dtype: bool	

```
movies[['color']].max()
Series([], dtype: float64)
movies.select_dtypes(['object']).fillna('')
                                                                  genres actor 1 name movie title actor 3 name
       color director name actor 2 name
                                                                                                                                              plot keywords
                                  Joel David Action|Adventure|Fantasy|Sci-Fi
    0 Color
                                                                           CCH Pounder
                                                                                               Avatar
                                                                                                           Wes Studi
                                                                                                                          avatar|future|marine|native|paraplegic
                    Cameron
                                                                                            Pirates of
                                                                                                                Jack
                                                                                                                          goddess|marriage ceremony|marriage http://ww
              Gore Verbinski Orlando Bloom
                                                 Action|Adventure|Fantasy
                                                                                          Caribbean:
                                                                                                           Davenport
                                                                                                                                                 proposal|pi.
                                                                                           At World's
                                                                                                End
                                                                               Christoph
                                                                                                           Stephanie
   2 Color
               Sam Mendes
                               Rory Kinnear
                                                  Action|Adventure|Thriller
                                                                                             Spectre
                                                                                                                            bomb|espionage|sequel|spy|terrorist
                                                                                   Waltz
                                                                                            The Dark
                                                                                                              Joseph deception|imprisonment|lawlessness|police http://ww
                               Christian Bale
   3 Color
                                                                              Tom Hardy
                                                                                              Knight
                                                            Action|Thriller
                                                                                               Rises
                                                                                           Star Wars:
                                                                                          Episode VII
                                                                                                                                                              http://w
                                                                            Doug Walker
                 Doug Walker
                                 Rob Walker
                                                            Documentary
                                                                                           The Force
                                                                                             Awakens
                                                                                              Signed
                                     Daphne
                 Scott Smith
 4911 Color
                                                                             Eric Mabius
                                                                                              Sealed
                                                                                                        Crystal Lowe
                                                                                                                             fraud|postal worker|prison|theft|trial
                                                          Comedy|Drama
                                     Zuniga
                                                                                            Delivered
                                                                                                               Sam
                                                                                                                                                              http://w
 4912 Color
                                Valorie Curry
                                              Crime|Drama|Mystery|Thriller
                                                                             Natalie Zea
                                                                                                                         cult|fbi|hideout|prison escape|serial killer
                                                                                                          Underwood
                                                                                            Following
                                                                                            A Plague
                   Benjamin
                                     Maxwell
                                                                                                               David
 4913 Color
                                                                           Eva Boehnke
                                                      Drama|Horror|Thriller
                     Roberds
                                                                                             Pleasant
                                                                                            Shanghai
                                                                                                         Eliza Coupe
 4914 Color
                  Daniel Hsia
                              Daniel Henney
                                                 Comedy|Drama|Romance
                                                                               Alan Ruck
                                                                                              Calling
                                                                                             My Date
                                                                                                                        actress name in title|crush|date|four word
                                                                             John August
 4915 Color
                                                                                                           Jon Gunn
                                                            Documentary
                                   Herzlinger
```

True

movies.isna().sum().sum()

2654

movies.isna().any() color True director name True num critic True duration True director fb True actor 3 fb True actor 2 name True actor 1 fb True True gross False genres actor 1 name True movie_title False num_voted_users False cast_total_fb False actor 3 name True facenumber in poster True plot_keywords True movie imdb link False True num user True language country True content_rating True budget True title year True actor 2 fb True imdb score False aspect ratio True movie fb False dtype: bool

movies.isna().any().any()

			()									
mov	ies[['co	lor']].max()									
Ser	ies([],	dtype: float	:64)								
		_			(11)							
mov	ies.	sele	ct_atypes([object']).fi	.lina(' ')							
	C	olor	director_name	actor_2_name	genres	actor_1_name	movie_title	actor_3_name	plot_keywords		_	
	0 C	Color	James Cameron	Joel David Moore	Action Adventure Fantasy Sci- Fi	CCH Pounder	Avatar	Wes Studi	avatar future marine native paraplegic	http://w	v.	
	1 C	Color	Gore Verbinski	Orlando Bloom	Action Adventure Fantasy	Johnny Depp	Pirates of the Caribbean: At World's End	Jack Davenport	goddess marriage ceremony marriage proposal pi		<pre>(movies.select_dty .fillna('') .max())</pre>	pes(['object'])
	2 C	Color	Sam Mendes	Rory Kinnear	Action Adventure Thriller	Christoph Waltz	Spectre	Stephanie Sigman	bomb espionage sequel spy terrorist		color director_name	Color Étienne Faure
	3 C	Color	Christopher Nolan	Christian Bale	Action Thriller	Tom Hardy	The Dark Knight Rises	Joseph Gordon-Levitt	deception imprisonment lawlessness police offi	пцр.	actor_2_name genres actor_1_name	Zubaida Sahar Western Óscar Jaenada
	4		Doug Walker	Rob Walker	Documentary	Doug Walker	Star Wars: Episode VII - The Force Awakens			http:/	movie_title actor_3_name plot_keywords	Æon Flux Óscar Jaenada zombie zombie spoof
											movie_imdb_link	http://www.imdb.com/title/tt5574490/?ref_=fn_t
491	1 C	Color	Scott Smith	Daphne Zuniga	Comedy Drama	Eric Mabius	Signed Sealed Delivered	Crystal Lowe	fraud postal worker prison theft trial	http:/	language country content_rating	Zulu West Germany X
491	2 C	Color		Valorie Curry	Crime Drama Mystery Thriller	Natalie Zea	The Following	Sam Underwood	cult fbi hideout prison escape serial killer	http:/	dtype: object	
491	3 C	Color	Benjamin Roberds	Maxwell Moody	Drama Horror Thriller	Eva Boehnke	A Plague So Pleasant	David Chandler		http://w	w	
491	4 C	olor	Daniel Hsia	Daniel Henney	Comedy Drama Romance	Alan Ruck	Shanghai Calling	Eliza Coupe		http://w	w	
491	5 C	olor	Jon Gunn	Brian Herzlinger	Documentary	John August	My Date with Drew	Jon Gunn	actress name in title crush date four word tit	http://w	w	18

```
colleges = pd.read csv('data/college.csv')
# colleges + 5
colleges = pd.read csv('data/college.csv', index col='INSTNM')
college ugds = colleges.filter(like='UGDS ')
college ugds.head()
                               UGDS_WHITE UGDS_BLACK UGDS_HISP UGDS_ASIAN UGDS_AIAN UGDS_NHPI UGDS_2MOR UGDS_NRA UGDS_UNKN
                      INSTNM
                                                                                                                  0.0000
                                                                                                                                           0.0138
       Alabama A & M University
                                     0.0333
                                                   0.9353
                                                               0.0055
                                                                             0.0019
                                                                                         0.0024
                                                                                                     0.0019
                                                                                                                              0.0059
        University of Alabama at
                                     0.5922
                                                   0.2600
                                                               0.0283
                                                                             0.0518
                                                                                         0.0022
                                                                                                     0.0007
                                                                                                                  0.0368
                                                                                                                              0.0179
                                                                                                                                           0.0100
                   Birmingham
             Amridge University
                                                                                                     0.0000
                                                                                                                  0.0000
                                                                                                                              0.0000
                                                                                                                                           0.2715
                                     0.2990
                                                   0.4192
                                                               0.0069
                                                                             0.0034
                                                                                         0.0000
        University of Alabama in
                                     0.6988
                                                   0.1255
                                                               0.0382
                                                                             0.0376
                                                                                         0.0143
                                                                                                     0.0002
                                                                                                                  0.0172
                                                                                                                              0.0332
                                                                                                                                           0.0350
                     Huntsville
        Alabama State University
                                                               0.0121
                                                                                                     0.0006
                                                                                                                  0.0098
                                                                                                                                           0.0137
                                     0.0158
                                                   0.9208
                                                                             0.0019
                                                                                         0.0010
                                                                                                                              0.0243
```

```
name = 'Northwest-Shoals Community College'
college_ugds.loc[name]
UGDS WHITE
              0.7912
UGDS BLACK
             0.1250
UGDS HISP
             0.0339
UGDS ASIAN
             0.0036
UGDS AIAN
             0.0088
UGDS NHPI
             0.0006
UGDS_2MOR
             0.0012
UGDS NRA
             0.0033
UGDS UNKN
             0.0324
Name: Northwest-Shoals Community College, dtype: float64
college_ugds.loc[name].round(2)
UGDS_WHITE
             0.79
UGDS BLACK
             0.12
UGDS HISP
             0.03
UGDS ASIAN
             0.00
UGDS AIAN
             0.01
UGDS NHPI
             0.00
UGDS 2MOR
             0.00
UGDS NRA
             0.00
UGDS UNKN
             0.03
Name: Northwest-Shoals Community College, dtype: float64
(college_ugds.loc[name] + .0001).round(2)
UGDS WHITE
             0.79
UGDS_BLACK
             0.13
UGDS HISP
             0.03
UGDS ASIAN
             0.00
UGDS AIAN
             0.01
UGDS NHPI
             0.00
UGDS_2MOR
             0.00
UGDS NRA
             0.00
UGDS UNKN
             0.03
```

Name: Northwest-Shoals Community College, dtype: float64

college_ugds + .00501										
		UGDS_WHITE	UGDS_BLACK	UGDS_HISP	UGDS_ASIAN	UGDS_AIAN	UGDS_NHPI	UGDS_2MOR	UGDS_NRA	UGDS_UNKN
INS	NM									
Alabama A & M Unive	sity	0.03831	0.94031	0.01051	0.00691	0.00741	0.00691	0.00501	0.01091	0.01881
University of Alaban Birming		0.59721	0.26501	0.03331	0.05681	0.00721	0.00571	0.04181	0.02291	0.01501
Amridge Unive	sity	0.30401	0.42421	0.01191	0.00841	0.00501	0.00501	0.00501	0.00501	0.27651
University of Alabam Hunts		0.70381	0.13051	0.04321	0.04261	0.01931	0.00521	0.02221	0.03821	0.04001
Alabama State Unive	sity	0.02081	0.92581	0.01711	0.00691	0.00601	0.00561	0.01481	0.02931	0.01871
SAE Institute of Technology Franc		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Rasmussen College - Over	and Park	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
National Personal Trai Institute of Cleve		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Bay Area Medical Academy - Jose Satellite Loca		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Excel Learning Center Antonio S		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

7535 rows × 9 columns

.045+.005

0.0499999999999999

(college_ugds + .00501) // .01

	UGDS_WHITE	UGDS_BLACK	UGDS_HISP	UGDS_ASIAN	UGDS_AIAN	UGDS_NHPI	UGDS_2MOR	UGDS_NRA	UGDS_UNKN
INSTNM									
Alabama A & M University	3.0	94.0	1.0	0.0	0.0	0.0	0.0	1.0	1.0
University of Alabama at Birmingham	59.0	26.0	3.0	5.0	0.0	0.0	4.0	2.0	1.0
Amridge University	30.0	42.0	1.0	0.0	0.0	0.0	0.0	0.0	27.0
University of Alabama in Huntsville	70.0	13.0	4.0	4.0	1.0	0.0	2.0	3.0	4.0
Alabama State University	2.0	92.0	1.0	0.0	0.0	0.0	1.0	2.0	1.0
SAE Institute of Technology San Francisco	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Rasmussen College - Overland Park	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
National Personal Training Institute of Cleveland	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Bay Area Medical Academy - San Jose Satellite Location	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Excel Learning Center-San Antonio South	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

21

college ugds op round = (college ugds + .00501) // .01 / 100 college ugds op round.head()

UGDS_WHITE UGDS_BLACK UGDS_HISP UGDS_ASIAN UGDS_AIAN UGDS_NHPI UGDS_2MOR UGDS_NRA UGDS_UNKN INSTNM 0.03 0.94 0.01 0.00 0.00 0.0 0.00 0.01 0.01 Alabama A & M University University of Alabama at 0.59 0.26 0.03 0.05 0.00 0.0 0.04 0.02 0.01 Birmingham Amridge University 0.30 0.42 0.01 0.00 0.00 0.0 0.00 0.00 0.27 University of Alabama in 0.70 0.13 0.04 0.04 0.01 0.0 0.02 0.03 0.04 Huntsville Alabama State University 0.02 0.02 0.92 0.01 0.00 0.00 0.0 0.01 0.01

college ugds round = (college ugds + .00001).round(2) college_ugds_round

Jose Satellite Location

	UGDS_WHITE	UGDS_BLACK	UGDS_HISP	UGDS_ASIAN	UGDS_AIAN	UGDS_NHPI	UGDS_2MOR	UGDS_NRA	UGDS_UNKN
INSTNM									
Alabama A & M University	0.03	0.94	0.01	0.00	0.00	0.0	0.00	0.01	0.01
University of Alabama at Birmingham		0.26	0.03	0.05	0.00	0.0	0.04	0.02	0.01
Amridge University	0.30	0.42	0.01	0.00	0.00	0.0	0.00	0.00	0.27
University of Alabama in Huntsville		0.13	0.04	0.04	0.01	0.0	0.02	0.03	0.04
Alabama State University	0.02	0.92	0.01	0.00	0.00	0.0	0.01	0.02	0.01
SAE Institute of Technology San Francisco		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Rasmussen College - Overland Park		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
National Personal Training Institute of Cleveland		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
Bay Area Medical Academy - San	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

```
college_ugds_op_round.equals(college_ugds_round)
```

True

True

比較缺失值

np.nan == np.nan

False

None == None

True

np.nan > 5

False

5 > np.nan

False

np.nan != 5

True

```
college = pd.read_csv('data/college.csv', index_col='INSTNM')
college_ugds = college.filter(like='UGDS_')

college_ugds == .0019
```

UGDS_WHITE UGDS_BLACK UGDS_HISP UGDS_ASIAN UGDS_AIAN UGDS_NHPI UGDS_2MOR UGDS_NRA UGDS_UNKN

in 3 in in									
Alabama A & M University	False	False	False	True	False	True	False	False	False
University of Alabama at Birmingham	False								
Amridge University	False								
University of Alabama in Huntsville	False								
Alabama State University	False	False	False	True	False	False	False	False	False
SAE Institute of Technology San Francisco	False								
Rasmussen College - Overland Park	False								
National Personal Training Institute of Cleveland	False								
Bay Area Medical Academy - San Jose Satellite Location	False								
Excel Learning Center-San Antonio South	False								

7535 rows × 9 columns

INSTNM

比較缺失值

```
college_self_compare = college_ugds == college_ugds
college_self_compare.head()
```

UGDS_WHITE UGDS_BLACK UGDS_HISP UGDS_ASIAN UGDS_AIAN UGDS_NHPI UGDS_2MOR UGDS_NRA UGDS_UNKN

INSTNM

| Alabama A & M University | True |
|--|------|------|------|------|------|------|------|------|------|
| University of Alabama at
Birmingham | True |
| Amridge University | True |
| University of Alabama in
Huntsville | True |
| Alabama State University | True |

college_self_compare.all()

UGDS_WHITE False UGDS_BLACK False UGDS_HISP False UGDS_ASIAN False UGDS AIAN False UGDS NHPI False UGDS_2MOR False UGDS NRA False UGDS UNKN False dtype: bool

乍看沒有遺失值,但透過all()可以發現其實每個欄位都有遺失值。

比較缺失值

from pandas.testing import assert frame equal

True

assert_frame_equal(college_ugds, college_ugds) is None

```
(college_ugds == np.nan).sum()
UGDS WHITE
UGDS BLACK
UGDS HISP
UGDS ASIAN
UGDS AIAN
UGDS NHPI
UGDS 2MOR
UGDS NRA
UGDS UNKN
dtype: int64
college ugds.isna().sum()
UGDS WHITE
              661
UGDS BLACK
              661
UGDS HISP
              661
UGDS ASIAN
              661
UGDS AIAN
              661
```

661

661

661

661

UGDS NHPI

UGDS 2MOR

UGDS_NRA

UGDS UNKN

dtype: int64

相同位置的 nan 會當成同樣的元素 college_ugds.equals(college_ugds) True 與 college_ugds == .0019 相同 college_ugds.eq(.0019) UGDS WHITE UGDS BLACK UGDS HISP UGDS ASIAN UGDS AIAN UGDS NHPI UGDS 2MOR UGDS NRA UGDS UNKN INSTNM Alabama A & M University False False False True False True False False False University of Alabama at False False False False False False False False False Birmingham **Amridge University** False False False False False False False False False University of Alabama in False False False False False False False False False Alabama State University False False False True False False False False False SAE Institute of Technology San False False False False False False False False False Rasmussen College - Overland False False False False False False False False False National Personal Training False False False False False False False False False Institute of Cleveland Bay Area Medical Academy - San False False False False False False False False False Jose Satellite Location Excel Learning Center-San False False False False False False False False False Antonio South 7535 rows × 9 columns

單元測試,可正確比較 nan

轉置 DataFrame 運算的方向

```
college = pd.read_csv('data/college.csv', index_col='INSTNM')
college_ugds = college.filter(like='UGDS_')
college_ugds.head()
```

UGDS_WHITE UGDS_BLACK UGDS_HISP UGDS_ASIAN UGDS_AIAN UGDS_NHPI UGDS_2MOR UGDS_NRA UGDS_UNKN

INSTNM

Alabama A & M University	0.0333	0.9353	0.0055	0.0019	0.0024	0.0019	0.0000	0.0059	0.0138
University of Alabama at Birmingham	0.5922	0.2600	0.0283	0.0518	0.0022	0.0007	0.0368	0.0179	0.0100
Amridge University	0.2990	0.4192	0.0069	0.0034	0.0000	0.0000	0.0000	0.0000	0.2715
University of Alabama in Huntsville	0.6988	0.1255	0.0382	0.0376	0.0143	0.0002	0.0172	0.0332	0.0350
Alabama State University	0.0158	0.9208	0.0121	0.0019	0.0010	0.0006	0.0098	0.0243	0.0137

college_ugds.count()

UGDS_WHITE 6874 UGDS BLACK 6874 UGDS_HISP 6874 UGDS_ASIAN 6874 UGDS AIAN 6874 UGDS NHPI 6874 UGDS_2MOR 6874 UGDS NRA 6874 UGDS UNKN 6874 dtype: int64

轉置 DataFrame 運算的方向

```
INSTNM
Alabama A & M University 9
University of Alabama at Birmingham 9
Amridge University 9
University of Alabama in Huntsville 9
Alabama State University 9
dtype: int64
```

college_ugds.sum(axis='columns').head()

INSTNM Alabama A & M University 1.0000 University of Alabama at Birmingham 0.9999 Amridge University 1.0000 University of Alabama in Huntsville 1.0000 Alabama State University 1.0000 dtype: float64

college_ugds.median(axis='index')

```
UGDS WHITE
              0.55570
UGDS BLACK
             0.10005
UGDS HISP
              0.07140
UGDS ASIAN
             0.01290
UGDS AIAN
             0.00260
UGDS NHPI
              0.00000
UGDS_2MOR
              0.01750
UGDS NRA
              0.00000
UGDS UNKN
              0.01430
dtype: float64
```

axis 參數預設為 0,即為 index。參數若為 1,則代表 columns。

```
college_ugds_cumsum = college_ugds.cumsum(axis=1)
college_ugds_cumsum.head()
```

	UGDS_WHITE	UGDS_BLACK	UGDS_HISP	UGDS_ASIAN	UGDS_AIAN	UGDS_NHPI	UGDS_2MOR	UGDS_NRA	UGDS_UNKN
INSTNM									
Alabama A & M University	0.0333	0.9686	0.9741	0.9760	0.9784	0.9803	0.9803	0.9862	1.0000
University of Alabama at Birmingham	0.5922	0.8522	0.8805	0.9323	0.9345	0.9352	0.9720	0.9899	0.9999
Amridge University	0.2990	0.7182	0.7251	0.7285	0.7285	0.7285	0.7285	0.7285	1.0000
University of Alabama in Huntsville	0.6988	0.8243	0.8625	0.9001	0.9144	0.9146	0.9318	0.9650	1.0000
Alabama State University	0.0158	0.9366	0.9487	0.9506	0.9516	0.9522	0.9620	0.9863	1.0000

案例演練:確定大學校園的多樣性

練習 & 回家作業