



# Medidas de Similitud o Similaridad



# ¿Similitud?

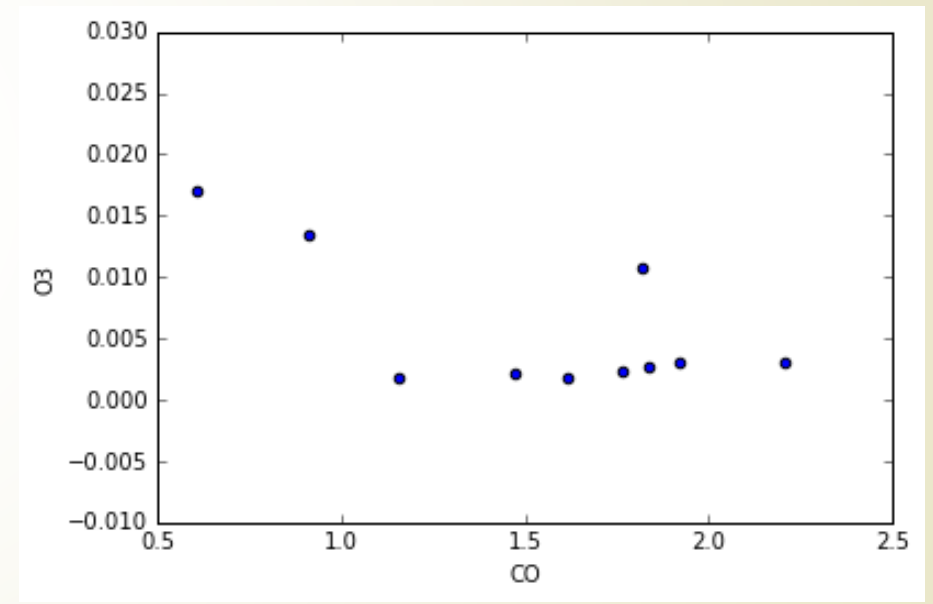
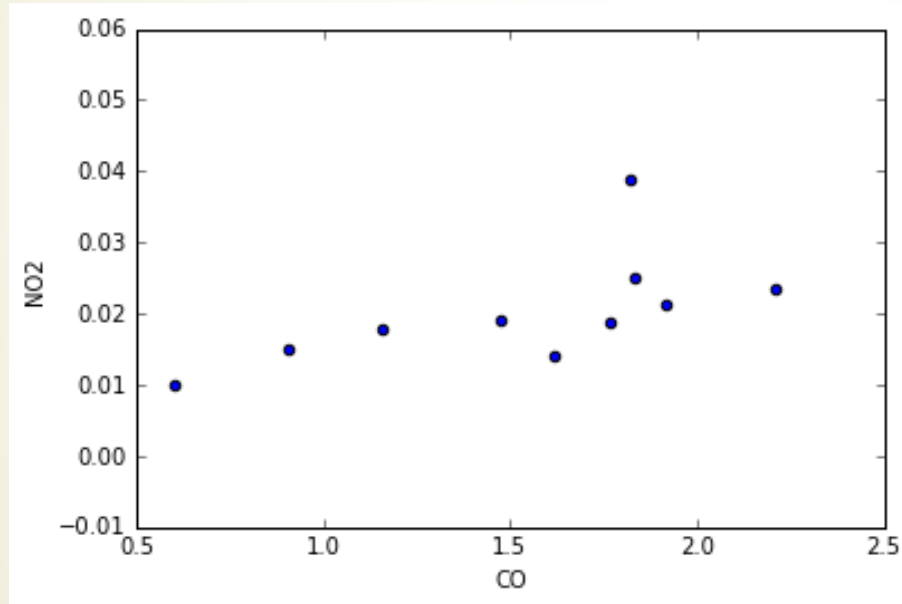
- ▶ Similitud se refiere a una forma de medir que tan parecido es un conjunto de datos con respecto a otro.
- ▶ El cálculo de la similitud depende del tipo de variable, y del significado o contexto de la misma.
- ▶ Las medidas de similitud con más frecuencia son utilizadas para variables con datos cuantitativos y datos binarios o de dos estados.
- ▶ Las variables categóricas pueden ser analizadas como variables binarias si se les aplica una transformación.

# Similitud para datos cuantitativos

CO	NO2	O3	PM10	SO2
0.603	0.01005	0.01695	43.95	0.00165
0.909	0.0151	0.01347	50.04	0.0019
1.835	0.02492	0.00273	46.23	0.00252
2.208	0.02338	0.00307	57.83	0.00298
1.473	0.01925	0.00223	56.66	0.00245
1.156	0.01782	0.00185	46.52	0.00182
1.766	0.01887	0.0024	40.54	0.00205
1.617	0.01397	0.00188	60.1	0.00208
1.917	0.02127	0.00308	57.22	0.00248
1.819	0.03882	0.01082	82.99	0.00228

- CO (**Monóxido de carbono**)
- NO2 (**Dióxido de Nitrógeno**)
- O3 (**Ozono a nivel del suelo**)
- PM10 (**Material Particulado 10 micrómetros**)
- SO2 (**Dióxido de azufre**)

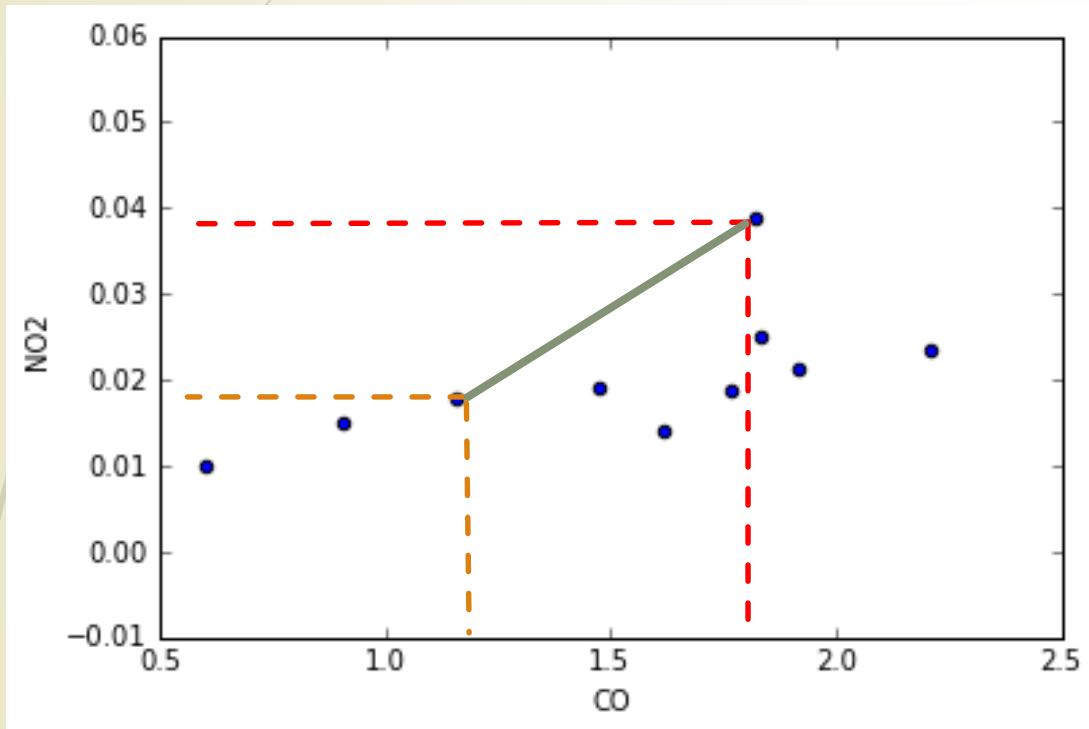
# Comparación mediante gráficos





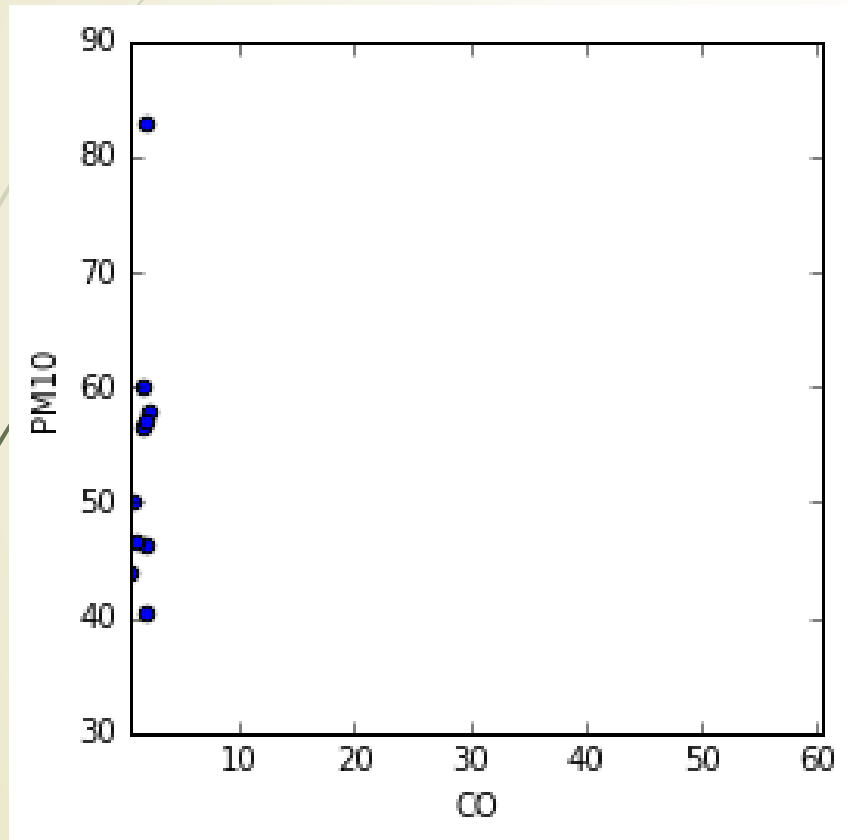
# Distancia Euclideana

# Distancia Euclideana



- La distancia Euclideana es la generalización del teorema de Pitágoras para más de 2 variables.
- $$d = \sqrt{(x_0 - x_1)^2 + (y_0 - y_1)^2}$$
- Esta medida de similitud puede ser usada para determinar cercanía entre dos muestras o dos variables.
- El problema principal se presenta cuando se tienen variables donde las escalas no son comparables o son muy diferentes.

# Distancia Euclideana



- Una solución es estandarizar los datos de la forma siguiente:

$$x^* = \frac{x - \mu_x}{\sigma_x}$$

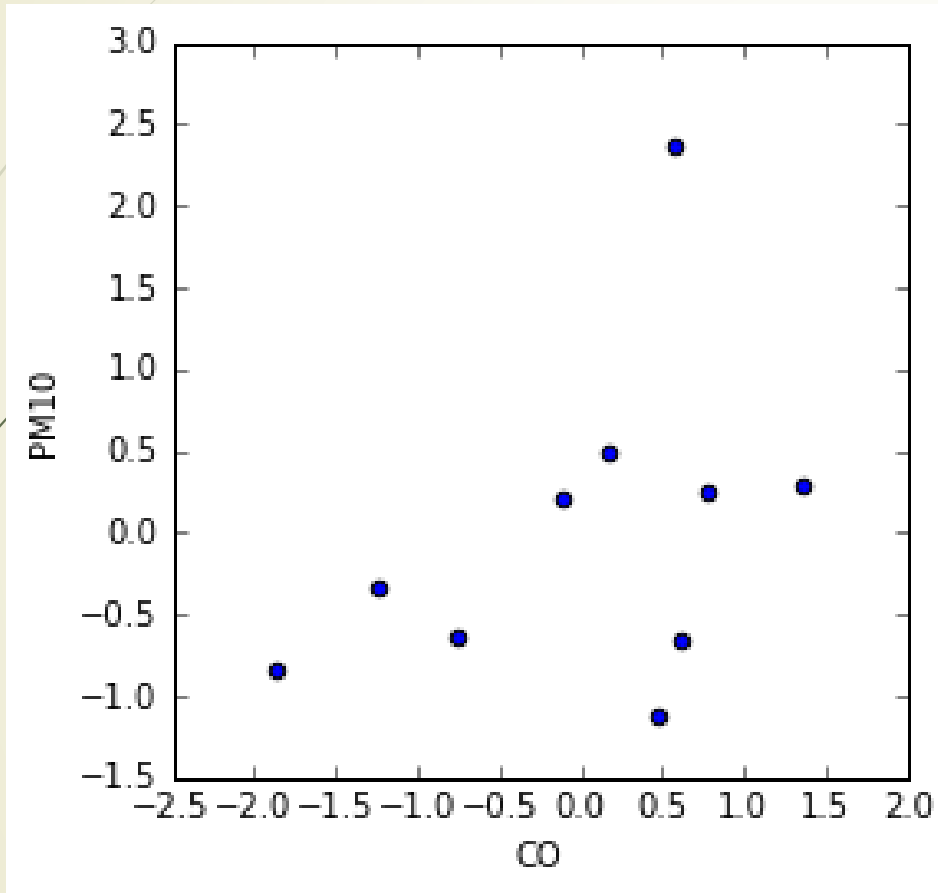
$$\mu = \begin{bmatrix} \mu_{CO} \\ \mu_{PM10} \end{bmatrix} = \begin{bmatrix} 1.5303 \\ 54.208 \end{bmatrix}$$

CO	PM10
0.603	43.95
0.909	50.04
1.835	46.23
2.208	57.83
1.473	56.66
1.156	46.52
1.766	40.54
1.617	60.1
1.917	57.22
1.819	82.99



CO	PM10
-1.85983221	-0.84590915
-1.24610563	-0.34370729
0.61111924	-0.65789269
1.35922386	0.29868229
-0.11492331	0.20220016
-0.75071195	-0.63397831
0.47272992	-1.12710921
0.1738892	0.48587412
0.77558192	0.24837964
0.57902896	2.37346043

# Distancia Euclideana



- Calculando la distancia Euclideana a los nuevos datos nos puede dar información que no tiene mucha dependencia de las escalas de los datos.
- $d^* = \sqrt{(x_0^* - x_1^*)^2 + (y_0^* - y_1^*)^2}$
- Ha esta distancia se le llama **distancia Euclideana Estandarizada**.



# Distancia Euclideana

- ▶ Una medida equivalente a la distancia Euclideana Estandarizada es la distancia Euclideana ponderada.

- ▶ 
$$d = \sqrt{\frac{1}{\sigma_x^2} (x_0 - x_1)^2 + \frac{1}{\sigma_y^2} (y_0 - y_1)^2}$$

# Matriz de Similitud

Matriz de similaridad de variables

	V0	V1	V2	V3	V4
V0		<b>0 2.59738013</b>	<b>5.47928451</b>	<b>3.29339382</b>	<b>1.5902661</b>
V1	<b>2.59738013</b>		<b>0 4.46595328</b>	<b>2.15890939</b>	<b>2.9584916</b>
V2	<b>5.47928451</b>	<b>4.46595328</b>		<b>0 4.06882391</b>	<b>5.21241236</b>
V3	<b>3.29339382</b>	<b>2.15890939</b>	<b>4.06882391</b>		<b>0 3.42203797</b>
V4	<b>1.5902661</b>	<b>2.9584916</b>	<b>5.21241236</b>	<b>3.42203797</b>	

Matriz de similaridad de muestras

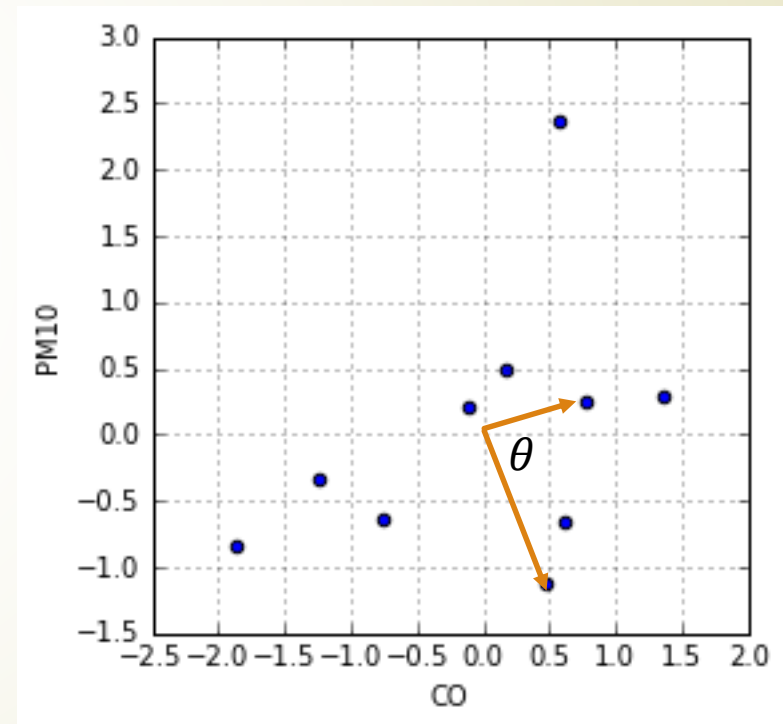
	S0	S1	S2	S3	S4	S5	S6	S7	S8	S9
S0		1.3469 0 1437	4.5620 8064	5.6323 2432	4.0400 3694	3.0923 247	3.7955 1369	3.7979 0485	4.5318 3541	5.7777 7361
S1	1.3469 1437		3.3373 7108	4.3640 639	2.7772 5954	2.1704 896	2.7876 6751	2.6700 2765	3.2480 1882	4.5769 0636
S2	4.5620 8064	3.3373 7108		1.6884 2376	1.3512 2355	2.4050 7474	1.4925 5027	2.1639 537	1.0382 0279	3.8374 0983
S3	5.6323 2432	4.3640 639	1.6884 2376		2.0619 1289	3.7884 2043	2.9351 1604	2.8326 6401	1.4111 71	3.7064 6943
S4	4.0400 3694	2.7772 5954	1.3512 2355	2.0619 1289		1.9087 6627	1.7675 9819	1.2173 4734	0.9430 8993	3.7258 612
S5	3.0923 247	2.1704 896	2.4050 7474	3.7884 2043	1.9087 6627		1.4494 1613	1.6658 4184	2.4686 7957	4.6688 576
S6	3.7955 1369	2.7876 6751	1.4925 5027	2.9351 1604	1.7675 9819	1.4494 1613		1.7588 7818	1.8047 4361	4.6103 4677
S7	3.7979 0485	2.6700 2765	2.1639 537	2.8326 6401	1.2173 4734	1.6658 4184	1.7588 7818		1.5279 0738	4.0583 3101
S8	4.5318 3541	3.2480 1882	1.0382 0279	1.4111 71	0.9430 8993	2.4686 7957	1.8047 4361	1.5279 0738		3.4151 7986
S9	5.7777 7361	4.5769 0636	3.8374 0983	3.7064 6943	3.7258 612	4.6688 576	4.6103 4677	4.0583 3101	3.4151 7986	



# Similaridad basada en el coseno

# Similaridad basado en coseno

- ▶ La similitud calculada en base al coseno tiene como fundamento el algebra vectorial.
- ▶ El producto punto esta definido como:  $a \cdot b = \|A\| \|B\| \cos\theta$ , donde A y B son vectores.
- ▶ La similaridad es calculada como:  
$$\text{similaridad} = \cos\theta = \frac{a \cdot b}{\|A\| \|B\|}$$
- ▶ La distancia es calculada como:  
$$d = 1 - \text{similaridad} = 1 - \frac{a \cdot b}{\|A\| \|B\|}$$



# Similaridad basado en coseno

Matriz de similaridad de variables

	V0	V1	V2	V3	V4
V0		0.37479909	1.66791993	0.60258016	0.14049701
V1	0.37479909		1.10804104	0.25893832	0.48625959
V2	1.66791993	1.10804104		0.91974044	1.50940237
V3	0.60258016	0.25893832	0.91974044		0.65057466
V4	0.14049701	0.48625959	1.50940237	0.65057466	

Matriz de similaridad de muestras

	S0	S1	S2	S3	S4	S5	S6	S7	S8	S9
S0		0.0124 1374	1.3161 7283	1.9779 1461	0.9101 5244	0.0374 2016	0.9702 7353	1.6965 2992	1.9931 709	1.6179 6351
S1	0.0124 1374		1.4612 6773	1.9986 0501	0.7548 2546	0.0919 3713	1.1276 5075	1.5751 7583	1.9991 6796	1.4867 9678
S2	1.3161 7283	1.4612 6773		0.4925 2674	1.9732 7197	1.0472 4513	0.0611 1635	1.4604 9589	0.5753 0269	1.5504 9237
S3	1.9779 1461	1.9986 0501	0.4925 2674		1.2960 2277	1.8846 8094	0.8201 5747	0.4688 2	0.0043 7937	0.5600 0548
S4	0.9101 5244	0.7548 2546	1.9732 7197	1.2960 2277		1.1834 1674	1.9928 4454	0.3479 5579	1.2054 3069	0.2724 9552
S5	0.0374 2016	0.0919 3713	1.0472 4513	1.8846 8094	1.1834 1674		0.7005 0742	1.8649 1442	1.9243 8932	1.8078 9987
S6	0.9702 7353	1.1276 5075	0.0611 1635	0.8201 5747	1.9928 4454	0.7005 0742		1.7379 1603	0.9129 0646	1.8042 2921
S7	1.6965 2992	1.5751 7583	1.4604 9589	0.4688 2	0.3479 5579	1.8649 1442	1.7379 1603		0.3919 3968	0.0054 4474
S8	1.9931 709	1.9991 6796	0.5753 0269	0.0043 7937	1.2054 3069	1.9243 8932	0.9129 0646	0.3919 3968		0.4779 8224
S9	1.6179 6351	1.4867 9678	1.5504 9237	0.5600 0548	0.2724 9552	1.8078 9987	1.8042 2921	0.0054 4474	0.4779 8224	

# Similaridad basado en coseno (dependencia de escala)

	V0	V1	V2	V3	V4
V0	0	0.37479909	1.66791993	0.60258016	0.14049701
V1	0.37479909	0	1.10804104	0.25893832	0.48625959
V2	1.66791993	1.10804104	0	0.91974044	1.50940237
V3	0.60258016	0.25893832	0.91974044	0	0.65057466
V4	0.14049701	0.48625959	1.50940237	0.65057466	0

	V0	V1	V2	V3	V4
V0	0	0.03942439	0.4293539	0.0410487	0.01556278
V1	0.03942439	0	0.33330297	0.02860085	0.04473988
V2	0.4293539	0.33330297	0	0.26793425	0.3313351
V3	0.0410487	0.02860085	0.26793425	0	0.02347511
V4	0.01556278	0.04473988	0.3313351	0.02347511	0



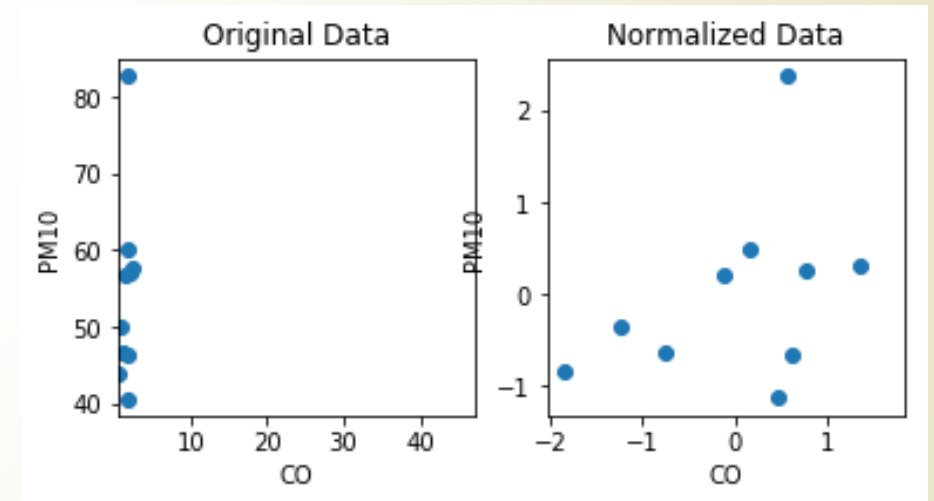
# Similaridad basado en coseno (dependencia de escala)

	S0	S1	S2	S3	S4	S5	S6	S7	S8	S9
S0		0.0124 1374	1.3161 7283	1.9779 1461	0.9101 5244	0.0374 2016	0.9702 7353	1.6965 2992	1.9931 709	1.6179 6351
S1	0.0124 1374		1.4612 6773	1.9986 0501	0.7548 2546	0.0919 3713	1.1276 5075	1.5751 7583	1.9991 6796	1.4867 9678
S2	1.3161 7283	1.4612 6773		0.4925 2674	1.9732 7197	1.0472 4513	0.0611 1635	1.4604 9589	0.5753 0269	1.5504 9237
S3	1.9779 1461	1.9986 0501	0.4925 2674		1.2960 2277	1.8846 8094	0.8201 5747	0.4688 2	0.0043 7937	0.5600 0548
S4	0.9101 5244	0.7548 2546	1.9732 7197	1.2960 2277		1.1834 1674	1.9928 4454	0.3479 5579	1.2054 3069	0.2724 9552
S5	0.0374 2016	0.0919 3713	1.0472 4513	1.8846 8094	1.1834 1674		0.7005 0742	1.8649 1442	1.9243 8932	1.8078 9987
S6	0.9702 7353	1.1276 5075	0.0611 1635	0.8201 5747	1.9928 4454	0.7005 0742		1.7379 1603	0.9129 0646	1.8042 2921
S7	1.6965 2992	1.5751 7583	1.4604 9589	0.4688 2	0.3479 5579	1.8649 1442	1.7379 1603		0.3919 3968	0.0054 4474
S8	1.9931 709	1.9991 6796	0.5753 0269	0.0043 7937	1.2054 3069	1.9243 8932	0.9129 0646	0.3919 3968		0.4779 8224
S9	1.6179 6351	1.4867 9678	1.5504 9237	0.5600 0548	0.2724 9552	1.8078 9987	1.8042 2921	0.0054 4474	0.4779 8224	0

	S0	S1	S2	S3	S4	S5	S6	S7	S8	S9
S0		9.8754 E-06	0.0003 3675	0.0002 9872	7.5301 E-05	6.1884 E-05	0.0004 4444	8.6847 E-05	0.0001 9543	3.3583 E-05
S1	9.8754 E-06		0.0002 313	0.0001 9997	3.0637 E-05	2.2317 E-05	0.0003 2182	3.8152 E-05	0.0001 1745	7.0362 E-06
S2	0.0003 3675	0.0002 313		1.1396 E-06	9.3579 E-05	0.0001 0993	7.459E -06	8.1578 E-05	1.911E -05	0.0001 5766
S3	0.0002 9872	0.0001 9997	1.1396 E-06		7.4066 E-05	8.8682 E-05	1.4429 E-05	6.3434 E-05	1.0917 E-05	0.0001 3199
S4	7.5301 E-05	3.0637 E-05	9.3579 E-05	7.4066 E-05		6.577E -07	0.0001 5388	4.1164 E-07	2.8113 E-05	8.309E -06
S5	6.1884 E-05	2.2317 E-05	0.0001 0993	8.8682 E-05	6.577E -07		0.0001 7465	2.11E- 06	3.7371 E-05	4.2913 E-06
S6	0.0004 4444	0.0003 2182	7.459E -06	1.4429 E-05	0.0001 5388	0.0001 7465		0.0001 3837	5.0447 E-05	0.0002 3369
S7	8.6847 E-05	3.8152 E-05	8.1578 E-05	6.3434 E-05	4.1164 E-07	2.11E- 06	0.0001 3837		2.1721 E-05	1.2419 E-05
S8	0.0001 9543	0.0001 1745	1.911E -05	1.0917 E-05	2.8113 E-05	3.7371 E-05	5.0447 E-05	2.1721 E-05		6.6989 E-05
S9	3.3583 E-05	7.0362 E-06	0.0001 5766	0.0001 3199	8.309E -06	4.2913 E-06	0.0002 3369	1.2419 E-05	6.6989 E-05	0

# Similaridad basado en coseno

- ▶ La similaridad tiene una dependencia en la escala de los datos.
- ▶ Pero esa dependencia solo afecta en la magnitud de los índices de similaridad, ya que el coseno depende del ángulo entre los puntos en el espacio.
- ▶ La proporción de los índices debe de guardar la misma proporción.







# Similaridad basada en correlación

# Correlación

- ▶ Existen varios tipos de coeficientes de correlación:
  - ▶ Coeficiente de correlación de Pearson
  - ▶ Coeficiente de correlación de Spearman
  - ▶ Correlación canónica
  - ▶ Coeficiente de Correlación Intraclass
  - ▶ Correlación de Kendall
  - ▶ Correlación de Jaspén
- ▶ El más conocido es el de Pearson

$$\rho_{x,y} = \frac{E[(x-\mu_x)(y-\mu_y)]}{\sigma_x\sigma_y}$$

- ▶ El coeficiente de correlación de Pearson no tiene dependencia de la escala de los datos, ya que implícitamente realiza la normalización de los datos.

# Correlación

Matriz de similitud de variables

	V0	V1	V2	V3	V4
V0		0.37479909	1.66791993	0.60258016	0.14049701
V1	0.37479909		1.10804104	0.25893832	0.48625959
V2	1.66791993	1.10804104		0.91974044	1.50940237
V3	0.60258016	0.25893832	0.91974044		0.65057466
V4	0.14049701	0.48625959	1.50940237	0.65057466	

Matriz de similitud de muestras

	S0	S1	S2	S3	S4	S5	S6	S7	S8	S9
S0		0.00478325	1.70752283	1.79085145	1.70269887	0.98354635	1.43780599	1.40425373	1.91199485	1.03725271
S1	0.00478325		1.73234964	1.81225543	1.66638952	0.94585996	1.51051933	1.41031466	1.93428097	0.96812732
S2	1.70752283	1.73234964		0.21260221	0.58375771	1.1356771	0.25212329	1.25079536	0.31479928	1.40082025
S3	1.79085145	1.81225543	0.21260221		0.21122927	1.579606	0.54172813	0.72862535	0.07208825	1.5442529
S4	1.70269887	1.66638952	0.58375771	0.21122927		1.45874345	1.1484746	0.59443504	0.26465274	1.13078433
S5	0.98354635	0.94585996	1.1356771	1.579606	1.45874345		1.00563954	1.27383457	1.39622728	0.1523876
S6	1.43780599	1.51051933	0.25212329	0.54172813	1.1484746	1.00563954		1.14290993	0.51116231	1.43717457
S7	1.40425373	1.41031466	1.25079536	0.72862535	0.59443504	1.27383457	1.14290993		0.47449784	0.91058237
S8	1.91199485	1.93428097	0.31479928	0.07208825	0.26465274	1.39622728	0.51116231	0.47449784		1.32300768
S9	1.03725271	0.96812732	1.40082025	1.5442529	1.13078433	0.1523876	1.43717457	0.91058237	1.32300768	



# Bases de Datos Heterogeneas

# Bases de datos heterogéneas.

- ▶ Una base de datos heterogénea es una combinación de varios tipos de datos en una sola base de datos.
- ▶ Para la obtención de medidas de similitud, se puede hacer uso de variables auxiliares (dummy) para las variables categóricas.
- ▶ Accident\_Severity y Day\_of\_Week

Longitude	Latitude	Number_of_Vehicles	Number_of_Casualties	Accident_Severity	Day_of_Week
-0.198465	51.505538	1	1	3	2
-0.178838	51.491836	1	1	3	2
-0.20559	51.51491	1	1	3	2
-0.208327	51.514952	1	1	3	3
-0.206022	51.496572	2	1	2	6
-0.19361	51.500788	2	1	3	5
-0.173519	51.495171	2	1	3	5
-0.163542	51.492497	2	1	3	1
-0.21198	51.513659	2	1	3	6
-0.199786	51.5159	2	2	3	3

# Variables Dummy

Accident_Severity
3
3
3
3
2
3
3
3
3
3



Severity_Low	Severity_Medium	Severity_High
0	0	1
0	0	1
0	0	1
0	0	1
0	1	0
0	0	1
0	0	1
0	0	1
0	0	1
0	0	1

# Variables Dummy

Longitude	Latitude	Number_o f_Vehicles	Number_o f_Casualti es	Severity 2	Severity 3	Sunday	Monday	Tuesday	Thursday	Friday
-0.198465	51.505538	1	1	0	1	0	1	0	0	0
-0.178838	51.491836	1	1	0	1	0	1	0	0	0
-0.20559	51.51491	1	1	0	1	0	1	0	0	0
-0.208327	51.514952	1	1	0	1	0	0	1	0	0
-0.206022	51.496572	2	1	1	0	0	0	0	0	1
-0.19361	51.500788	2	1	0	1	0	0	0	1	0
-0.173519	51.495171	2	1	0	1	0	0	0	1	0
-0.163542	51.492497	2	1	0	1	1	0	0	0	0
-0.21198	51.513659	2	1	0	1	0	0	0	0	1
-0.199786	51.5159	2	2	0	1	0	0	1	0	0



# Variables Dummy

Longitude	Latitude	Number_of_Vehicles	Number_of_Casualties	Severity 2	Severity 3	Sunday	Monday	Tuesday	Thursday	Friday
-0.27323522	0.13559719	-1.161895	-0.31622777	-0.31622777	0.31622777	-0.31622777	1.44913767	-0.47434165	-0.47434165	-0.47434165
0.91926386	-1.23487764	-1.161895	-0.31622777	-0.31622777	0.31622777	-0.31622777	1.44913767	-0.47434165	-0.47434165	-0.47434165
-0.70613662	1.07298517	-1.161895	-0.31622777	-0.31622777	0.31622777	-0.31622777	1.44913767	-0.47434165	-0.47434165	-0.47434165
-0.87243152	1.07718602	-1.161895	-0.31622777	-0.31622777	0.31622777	-0.31622777	-0.621059	1.8973666	-0.47434165	-0.47434165
-0.73238412	-0.76118265	0.77459667	-0.31622777	2.84604989	-2.84604989	-0.31622777	-0.621059	-0.47434165	-0.47434165	1.8973666
0.02174532	-0.33949808	0.77459667	-0.31622777	-0.31622777	0.31622777	-0.31622777	-0.621059	-0.47434165	1.8973666	-0.47434165
1.24243615	-0.90131075	0.77459667	-0.31622777	-0.31622777	0.31622777	-0.31622777	-0.621059	-0.47434165	1.8973666	-0.47434165
1.84861964	-1.16876438	0.77459667	-0.31622777	-0.31622777	0.31622777	2.84604989	-0.621059	-0.47434165	-0.47434165	-0.47434165
-1.09438083	0.94786008	0.77459667	-0.31622777	-0.31622777	0.31622777	-0.31622777	-0.621059	-0.47434165	-0.47434165	1.8973666
-0.35349666	1.17200503	0.77459667	2.84604989	-0.31622777	0.31622777	-0.31622777	-0.621059	1.8973666	-0.47434165	-0.47434165





# Variables Dummy

- ▶ Coeficiente de similaridad(disimilaridad) general de Gower.
- ▶ Estandarizar cada variable y multiplicar todas las columnas correspondientes a variables dummy por el factor  $1/\sqrt{2}=0.7071$
- ▶ Este factor compensa la codificación 0/1.