



ESCUELA DE INGENIERÍA
FACULTAD DE INGENIERÍA

EDUCACIÓN
PROFESIONAL

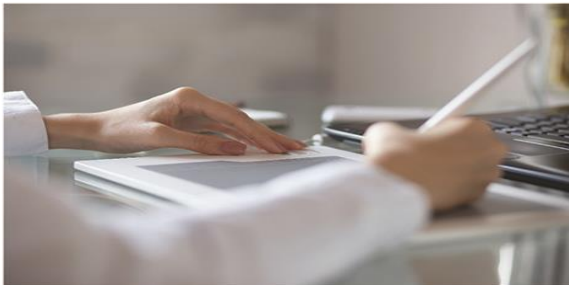
Diplomado en Big Data y Ciencias de Datos

Minería de Datos

Otros temas de su interés

Educación Profesional - Escuela de Ingeniería UC

Sebastián Raveau



Otros temas de su interés

Reducción de Dimensionalidad

Random Forests

Clustering Jerárquico

Minería de Texto

Procesamiento de Imágenes

Reducción de Dimensionalidad

Reducción de dimensionalidad

Supongamos que tenemos m variables explicativas

El objetivo es construir n nuevas variables (con $n < m$) que sean capten de la mejor manera posible la información contenida en las m variables originales

Un modelo que considere las n nuevas variables será más sencillo y tendrá un rendimiento similar a un modelo que considere las m variables originales

Análisis de Componentes Principales

El objetivo es explicar tanta variabilidad de las variables originales como sea posible

Mientras más componentes tengamos, más variabilidad explicaremos

Las componentes son una combinación lineal de las variables originales

Las componentes generadas son independientes, lo que facilita su inclusión en los modelos

Ejemplo

Encuesta de satisfacción a viajeros de una aerolínea

10 preguntas con nota de 1 a 7 respecto de:

Puntualidad

Entretenimiento a bordo

Comodidad de asientos

Programa de viajero frecuente

Comida

Frecuencia de vuelos

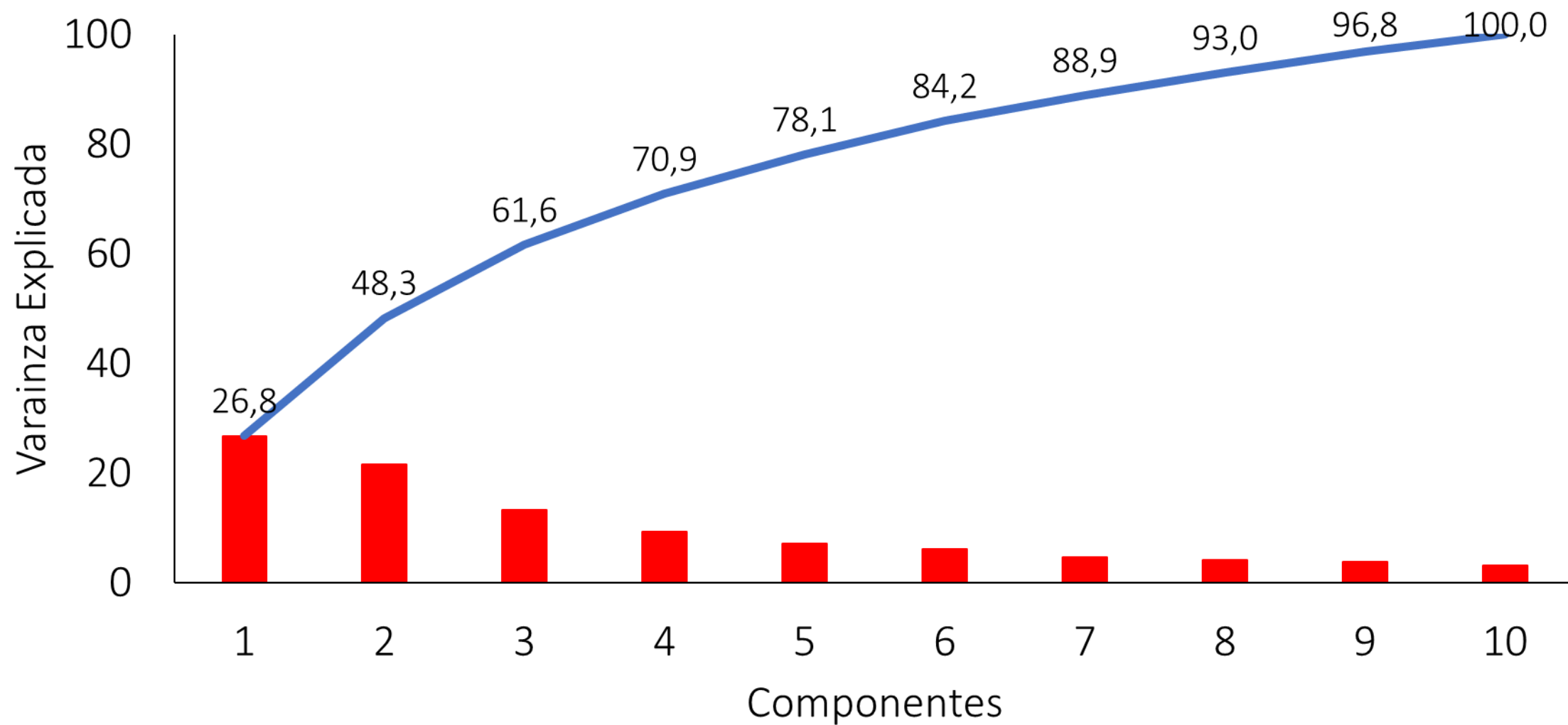
Simpatía de los tripulantes

Facilidad de maletas

Precio

Facilidad de pago

Ejemplo



Ejemplo

Variable	Componentes			
	1	2	3	4
Puntualidad	0,954	-0,004	0,041	0,153
Comodidad Asientos	0,037	0,965	0,001	-0,114
Comida	0,112	0,935	0,141	0,041
Simpatía Tripulación	-0,062	0,031	-0,121	0,133
Precio	0,156	-0,149	0,934	0,012
Entretenimiento a Bordo	-0,259	0,948	-0,041	0,121
Programa Viajero Frecuente	-0,131	0,220	0,112	0,934
Frecuencia de Vuelos	0,921	0,123	0,231	-0,032
Facilidad de Maletas	0,111	-0,031	0,099	0,966
Facilidad de Pago	0,011	0,044	-0,118	0,034

Random Forests

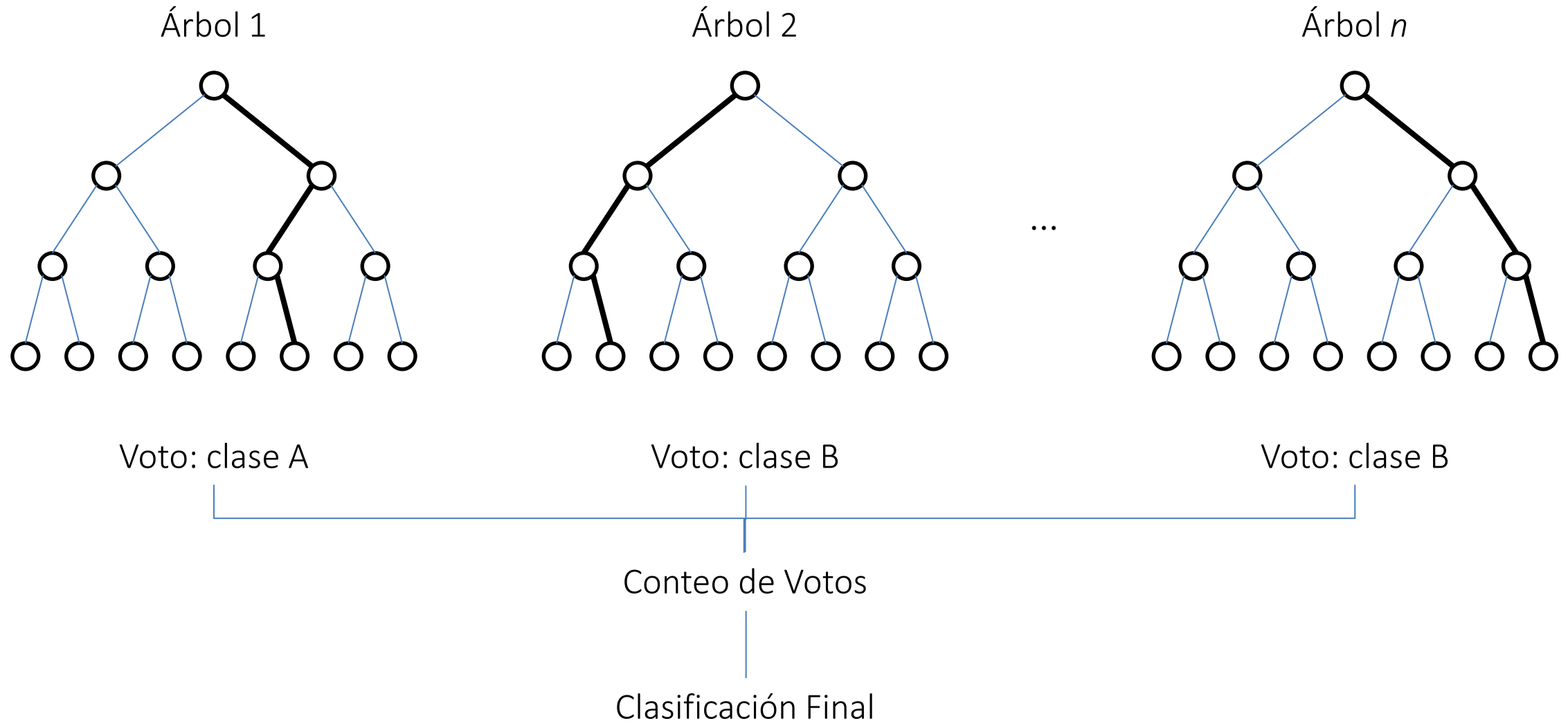
Random Forests

Un bosque está compuesto de varios árboles

En Random Forest se ejecutan varios algoritmos de árboles de decisión (clasificación o regresión) en lugar de sólo uno

Los árboles difieren en especificación y cada uno aporta un voto para la decisión final

Random Forests



Ventajas

Tiende a evitar sobreajuste

Decisiones más precisas y mejor informadas

Funciona bien para grandes volúmenes de datos y gran cantidad de variables

Permite identificar variables más significativas

Desventajas

Resultados más difíciles de interpretar (no hay un único árbol)

Computacionalmente costoso

Sobreajuste de datos atípicos

Clustering Jerárquico

Clustering Jerárquico

Otro método de clustering, alternativo a K Means

Existen dos tipos:

Aglomerativo: las observaciones se van juntando de una a la vez

Divisivo: las observaciones se van separando de una a la vez

Clustering Jerárquico Aglomerativo

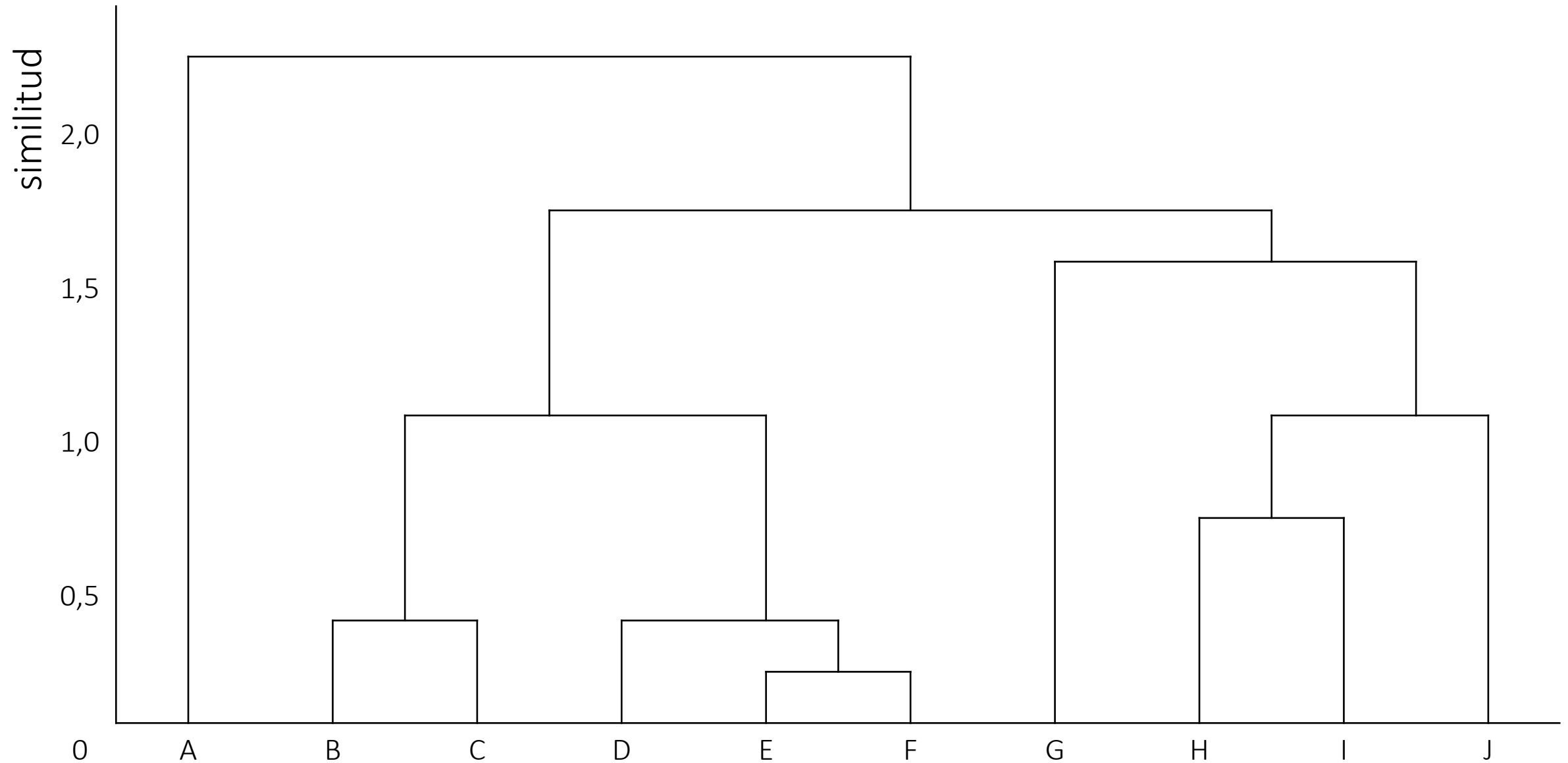
Inicialmente cada observación es su propio clúster

En cada iteración unimos dos clústeres de acuerdo a una métrica de similitud

Tras cada iteración actualizamos la métrica de similitud

El algoritmo termina cuando sólo queda un clúster con todas las observaciones

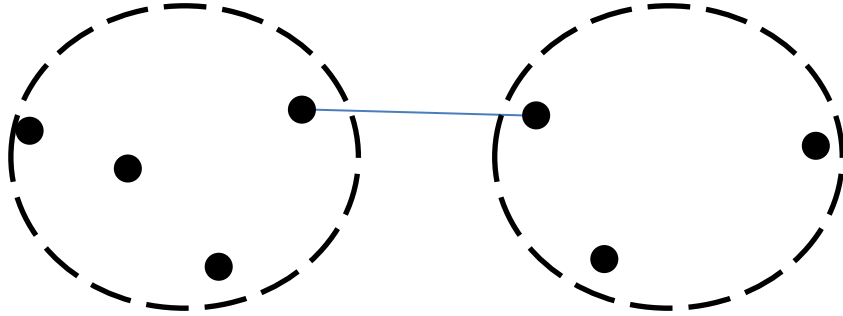
Dendrograma



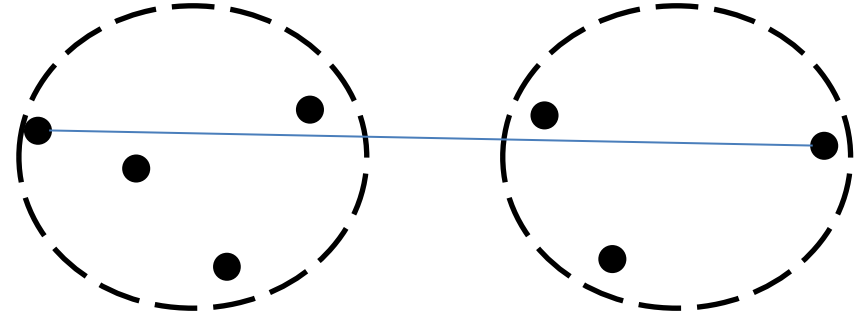
Métricas de similitud

Métricas de distancia

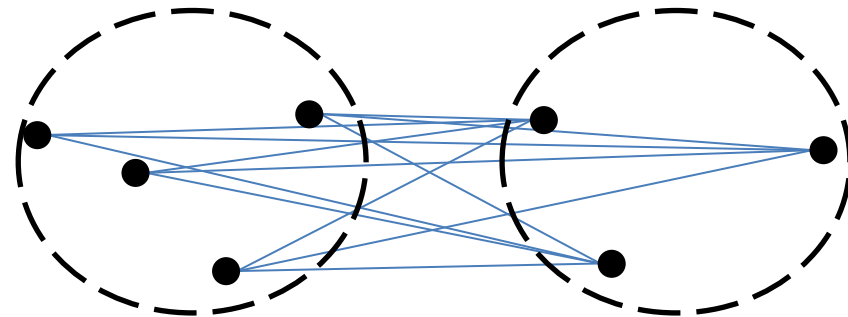
Mínimo



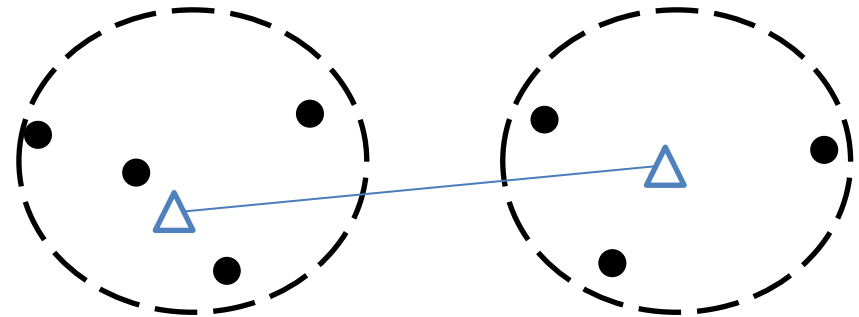
Máximo



Promedio



Centroides



Otras posibles métricas (correlación, desviación estándar, etc.)

Minería de Texto

Minería de Texto

Las bases de datos con las que hemos trabajado se ven así:

Obs.	Var. 1	Var. 2	Var. 3	Var. 4	Var. 5	...
1	6	1	8	1	5	...
2	5	3	2	4	8	...
3	7	6	1	3	7	...
4	3	0	4	9	6	...
5	1	2	6	0	3	...
6	9	1	2	3	4	...
7	4	4	2	1	5	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮

Minería de Texto

Pero los documentos se ven así

En un lugar de la Mancha, de cuyo nombre no quiero acordarme, no ha mucho tiempo que vivía un hidalgo de los de lanza en astillero, adarga antigua, rocín flaco y galgo corredor. Una olla de algo más vaca que carnero, salpicón las más noches, duelos y quebrantos los sábados, lantejas los viernes, algún palomino de añadidura los domingos, consumían las tres cuartas partes de su hacienda. El resto della concluían sayo de velarte, calzas de velludo para las fiestas, con sus pantuflos de lo mismo, y los días de entresemana se honraba con su vellorí de lo más fino. Tenía en su casa una ama que pasaba de los cuarenta, y una sobrina que no llegaba a los veinte, y un mozo de campo y plaza, que así ensillaba el rocín como tomaba la podadera. Frisaba la edad de nuestro hidalgo con los cincuenta años; era de complexión recia, seco de carnes, enjuto de rostro, gran madrugador y amigo de la caza. Quieren decir que tenía el sobrenombre de Quijada, o Quesada, que en esto hay alguna diferencia en los autores que deste caso escriben; aunque, por conjeturas verosímiles, se deja entender que se llamaba Quejana. Pero esto importa poco a nuestro cuento; basta que en la narración dél no se salga un punto de la verdad.

¿Cómo representar documentos en forma numérica?

Corpus

Un corpus es un conjunto de documentos

Ejemplo:

Documento 1: “un auto rojo”

Documento 2: “un tomate rojo y un globo rojo”

Documento 3: “un plátano amarillo y un tomate verde”

Vocabulario

Un vocabulario es una secuencia ordenada de palabras con un identificador único

En nuestro ejemplo:

ID	Palabra
1	amarillo
2	auto
3	globo
4	plátano
5	rojo
6	tomate
7	un
8	verde
9	y

Bolsa de palabras

Podemos representar cada documento como una bolsa de palabras en forma matricial

										ID	Palabra
Documento 1:	“un auto rojo”									1	amarillo
Documento 2:	“un tomate rojo y un globo rojo”									2	auto
Documento 3:	“un plátano amarillo y un tomate verde”									3	globo
										4	plátano
										5	rojo
										6	tomate
										7	un
Doc.	1	2	3	4	5	6	7	8	9	8	verde
1	0	1	0	0	1	0	1	0	0	9	y
2	0	0	1	0	2	1	2	0	1		
3	1	0	0	1	0	1	2	1	1		

Minería de Texto

¿Qué podemos hacer con estos datos?

Aprendizaje supervisado (e.g. análisis de sentimientos)

Clasificación

Análisis no supervisado

Clustering

Minería de Texto

Limpieza y preprocesamiento de los datos

Mayúsculas y minúsculas

Pronombres, preposiciones y otras palabras vacías

Acentos y puntuación

Stemming

Procesamiento de Imágenes

Procesamiento de Imágenes


Así como la unidad base de un texto son sus palabras, la unidad base de una imagen son sus píxeles

Cada píxel tiene básicamente dos tipos de información:


- Información espacial

- Información cromática


Color




R = 73
G = 86
B = 102



R = 124
G = 132
B = 149

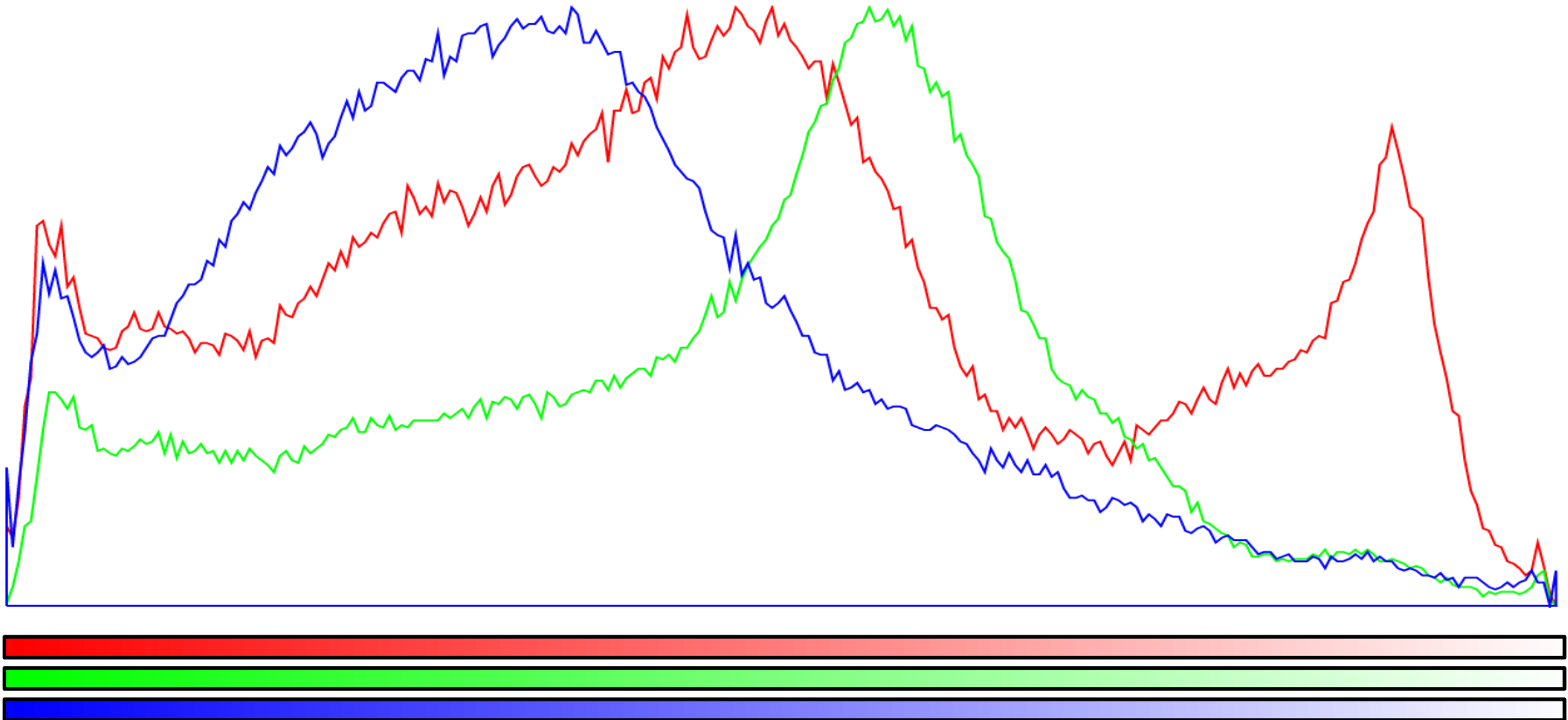


R = 223
G = 149
B = 112



R = 233
G = 238
B = 242

Histograma de colores



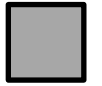
Intensidad



R = 84
G = 84
B = 84



R = 125
G = 125
B = 125

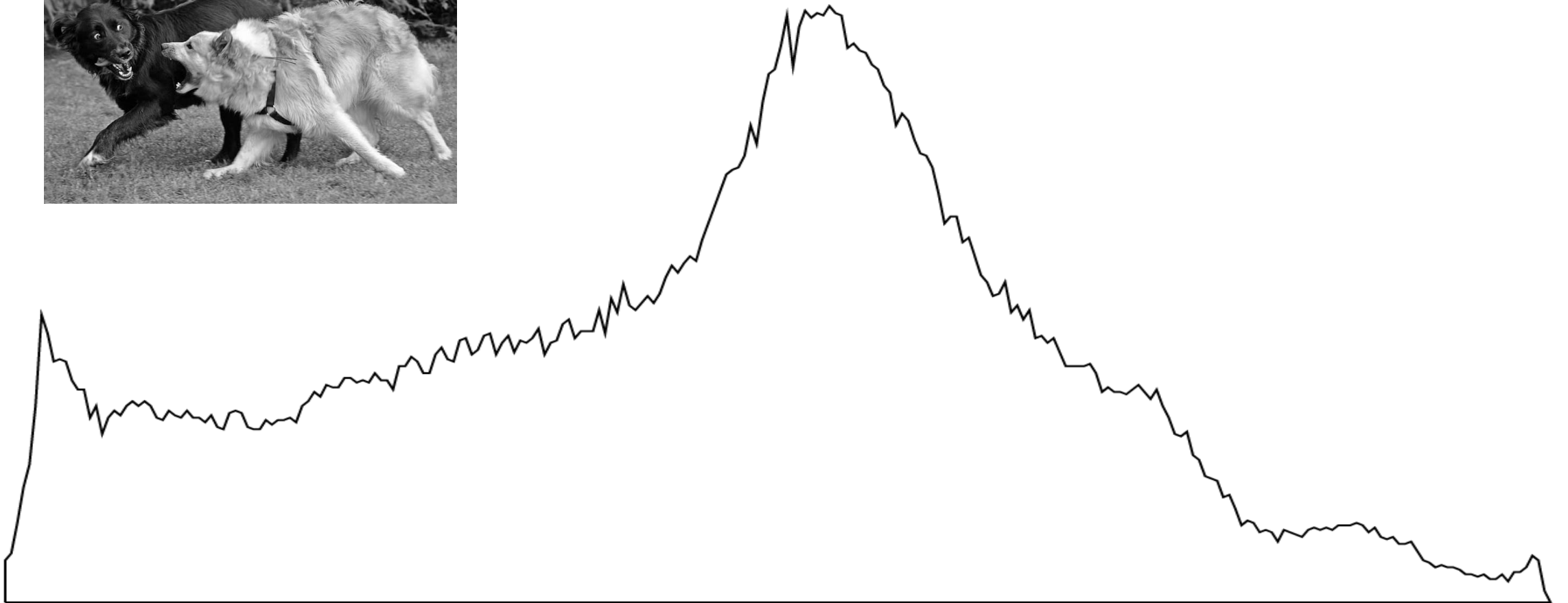


R = 167
G = 167
B = 167



R = 237
G = 237
B = 237

Histograma de intensidad



Procesamiento de Imágenes

Otros indicadores de textura basados en intensidad:

Bordes

Segundo momento angular

Contraste

Correlación

Entropía

Ejemplo: Clustering

$K = 2$



$K = 3$



$K = 4$



$K = 5$

