

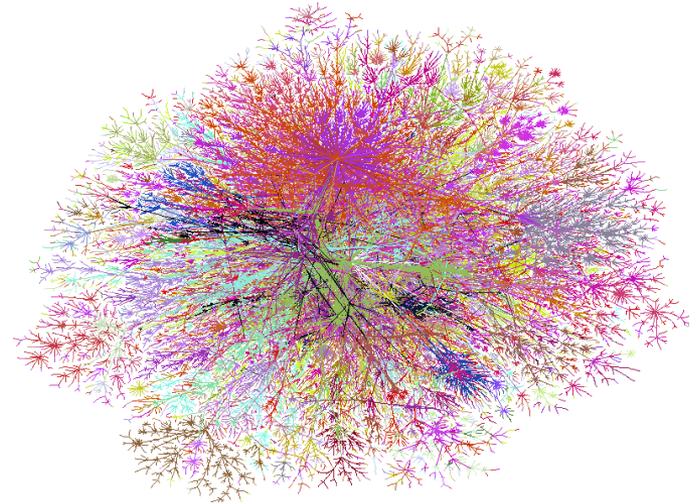
Web Data Representation

Web Graph, Text, Images, Metadata, Search spaces

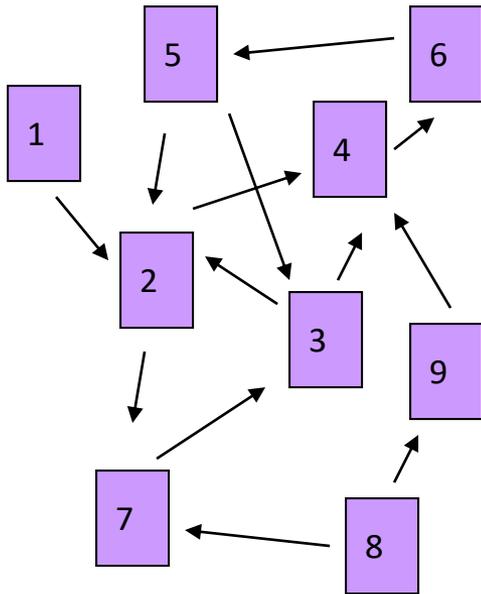
Web Search

The Web corpus

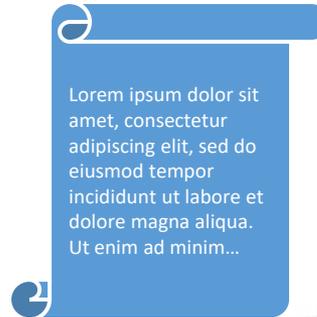
- No design/coordination
- Distributed content creation, linking, democratization of publishing
- Content includes truth, lies, obsolete information, contradictions ...
- Unstructured (text, html, ...), semi-structured (XML, annotated photos), structured (Databases)...
- Scale much larger than previous text corpora... but corporate records are catching up.
- Content can be dynamically generated



Web data



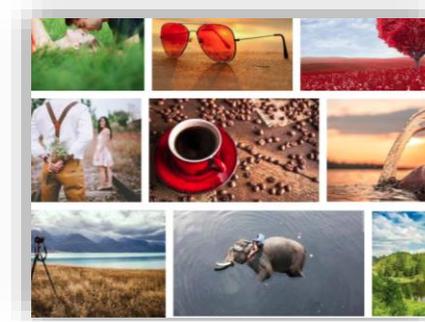
Links



Text



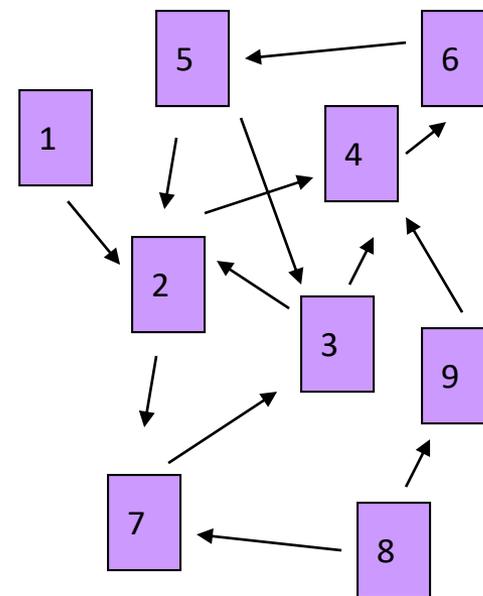
Ratings and clicks



Images/videos

The Web graph

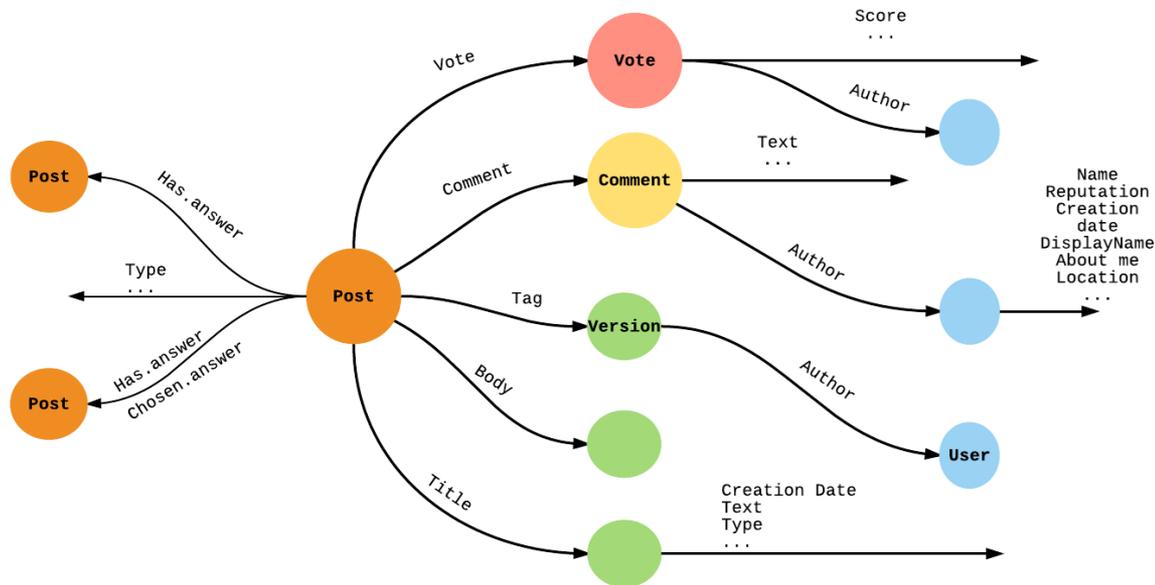
- Generally, the links can be explicit or computed by some function.
- The links can also be weighted by the similarity between pages (i.e. graph nodes in this case)
- Graphs are generally represented as a sparse matrix.
- There are many problems that make use of graph representations:
 - page importance, recommendation, reputation analysis...



1		1		1				
				1		1		
		1						1
					1			
			1					
	1						1	
							1	

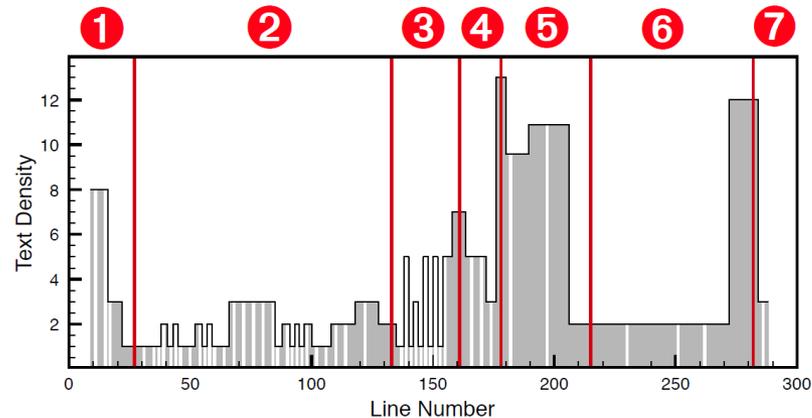
Graphs on the Web

- There are many types of graphs, besides hyperlinks.
- A single web page can be the source of many different information graphs:
 - Followers/followees, graph of named entities, authors, reply-to, ...



Web pages

- Web pages are divided into different parts (title, abstract, body, etc)
- Each part has a specific relevance to the main content
- A Web page can be divided by its HTML structure (e.g., <div> tags) or by its visual aspect.



Web page segmentation methods

- Segmenting visually
 - Cai, D., Yu, S., Wen, J. R., & Ma, W. Y. (2003). VIPS: A vision-based page segmentation algorithm.
- Linguistic approach
 - Kohlschütter, C. , Fankhauser, P., and Nejd, W. (2010). Boilerplate detection using shallow text features. ACM Web Search and Data Mining.
- Densitometric approach
 - Kohlschütter, C., and Nejd, W., (2008). A densitometric approach to web page segmentation. ACM Conference on Information and Knowledge Management (CIKM '08).

Text data

- Instead of aiming at fully understanding a text document, traditional search engines take a pragmatic approach and look at the most elementary text patterns
 - e.g. a simple histogram of words, also known as “bag-of-words”.
- Heuristics capture specific text patterns to improve search effectiveness
 - Enhances the simplicity of word histograms
- The most simple heuristics are stop-words removal and stemming

Character processing and stop-words

- Term delimitation
- Punctuation removal
- Numbers/dates
- Stop-words: remove words that are present in all documents
 - *a, and, are, as, at, be, but, by, for, if, in, into, is, it, no, not, of, on, or, such, that, the, their, then, there, these, they, this, to, was, will...*

Stemming and lemmatization

- Stemming: Reduce terms to their “roots” before indexing
 - “Stemming” suggest crude affix chopping
 - **e.g., automate(s), automatic, automation all reduced to automat.**
 - <http://tartarus.org/~martin/PorterStemmer/>
 - <http://snowball.tartarus.org/demo.php>
- Lemmatization: Reduce inflectional/variant forms to base form, e.g.,
 - *am, are, is* → *be*
 - *car, cars, car's, cars'* → *car*

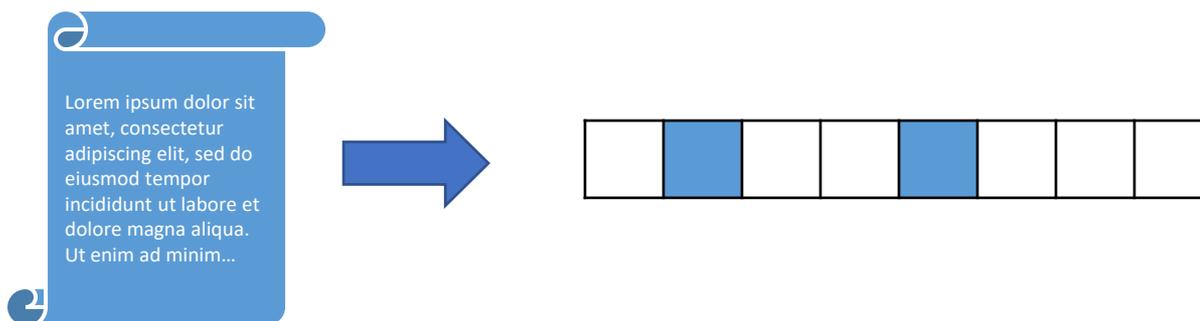
N-grams

- An n-gram is a sequence of items, e.g. characters, syllables or words.
- Can be applied to text spelling correction
 - *“interactive meida” >>>> “interactive media”*
- Can also be used as indexing tokens to improve Web page search
 - You can order the Google n-grams (6DVDs):
 - <http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>
- N-grams were under some criticism in NLP because they can add noise to information extraction tasks
 - ...but are widely successful in IR to infer document topics.

“Bag of Words” representation

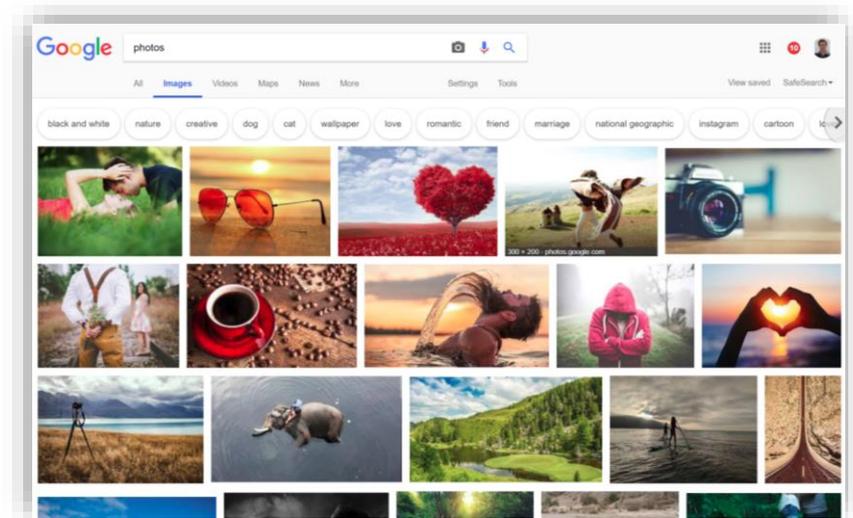
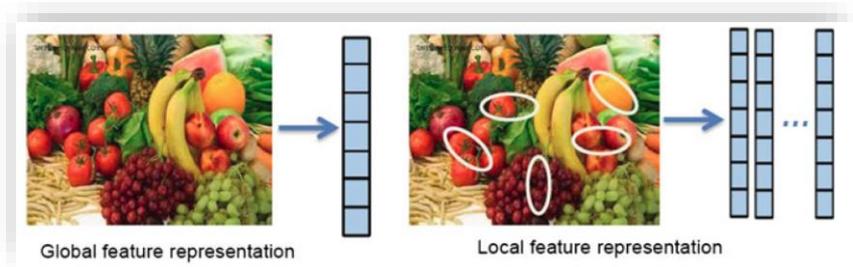
- After the text analysis steps, a document (e.g. Web page) is represented as a vector of terms and n-grams.
 - More complex low-level representations can be used

$$d = (w_1, \dots, w_L, ng_1, \dots, ng_M)$$



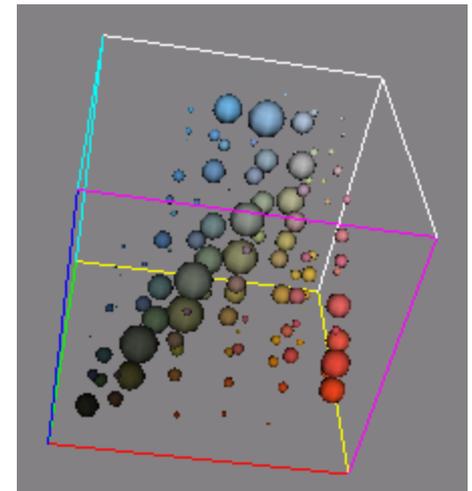
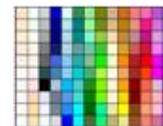
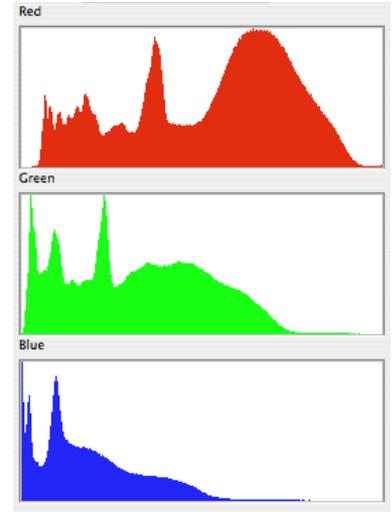
Visual data

- Visual information also needs to be processed and analysed.
- A compact representation of the image/video content is computed from it.
- This compact representation is then used to accomplish several tasks, e.g. search, categorization.



Histograms of colors

- Marginal color histograms consider color channels independently
 - The number of bins define the dimensionality of the space
- 3D colour histograms divide the space into small 3D boxes
 - The numbers of bins per dimension define the number of 3d bins



Color moments

- Color moments measure the statistical properties of the histogram:
 - Mean and variance (1st and 2nd moments)

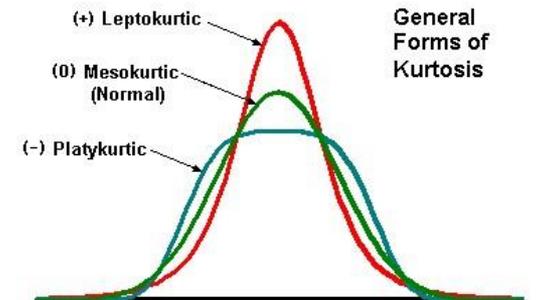
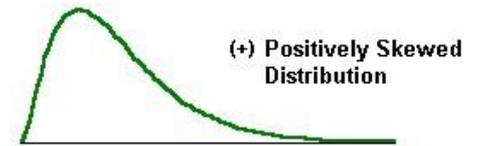
$$m_r = \sum \frac{(X_i - \bar{X})^r}{n}$$

- Skewness (3rd moment)

$$\text{Skewness} = \left[\frac{\sqrt{n(n-1)}}{n-2} \right] \times \frac{m_3}{m_2^{3/2}}$$

- Kurtosis (4th moment)

$$\text{Kurtosis} = \left[\frac{(n-1)(n+1)}{(n-2)(n-3)} \right] \times \frac{m_4}{(m_2)^2} - 3 \left[\frac{(n-1)^2}{(n-2)(n-3)} \right]$$

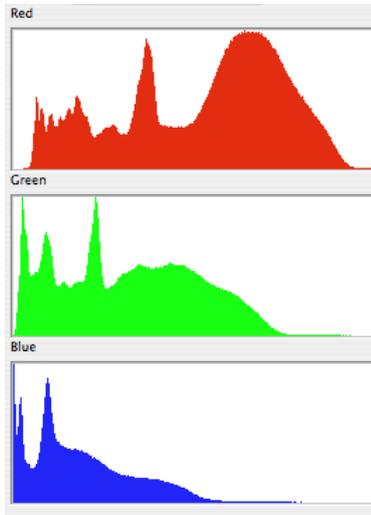


Example



Color moments

Marginal color histograms



$$\rightarrow d_{hR} = (bin_1, bin_2, \dots, bin_{16})$$

$$\rightarrow d_{hG} = (bin_1, bin_2, \dots, bin_{16})$$

$$\rightarrow d_{hB} = (bin_1, bin_2, \dots, bin_{16})$$

$$d_{cm} = (m_R, s_R^2, m_G, s_G^2, m_B, s_B^2)$$

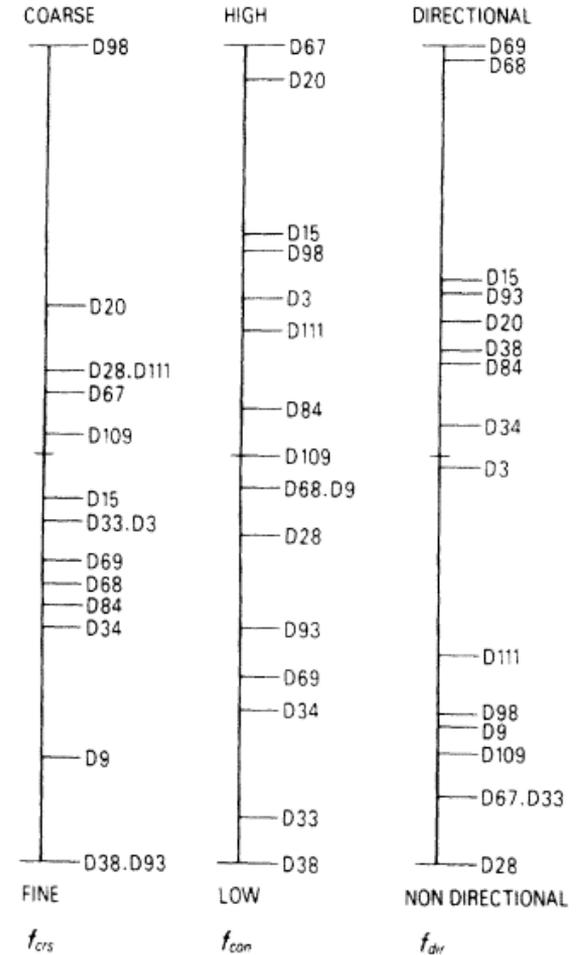
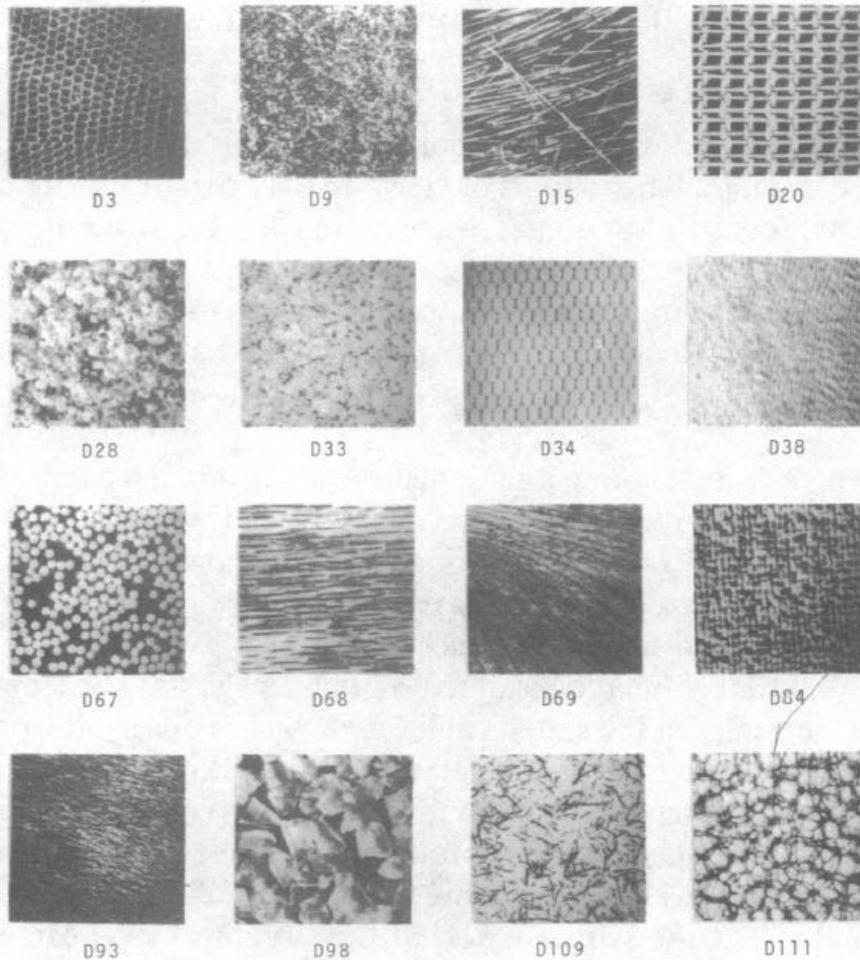
Textures



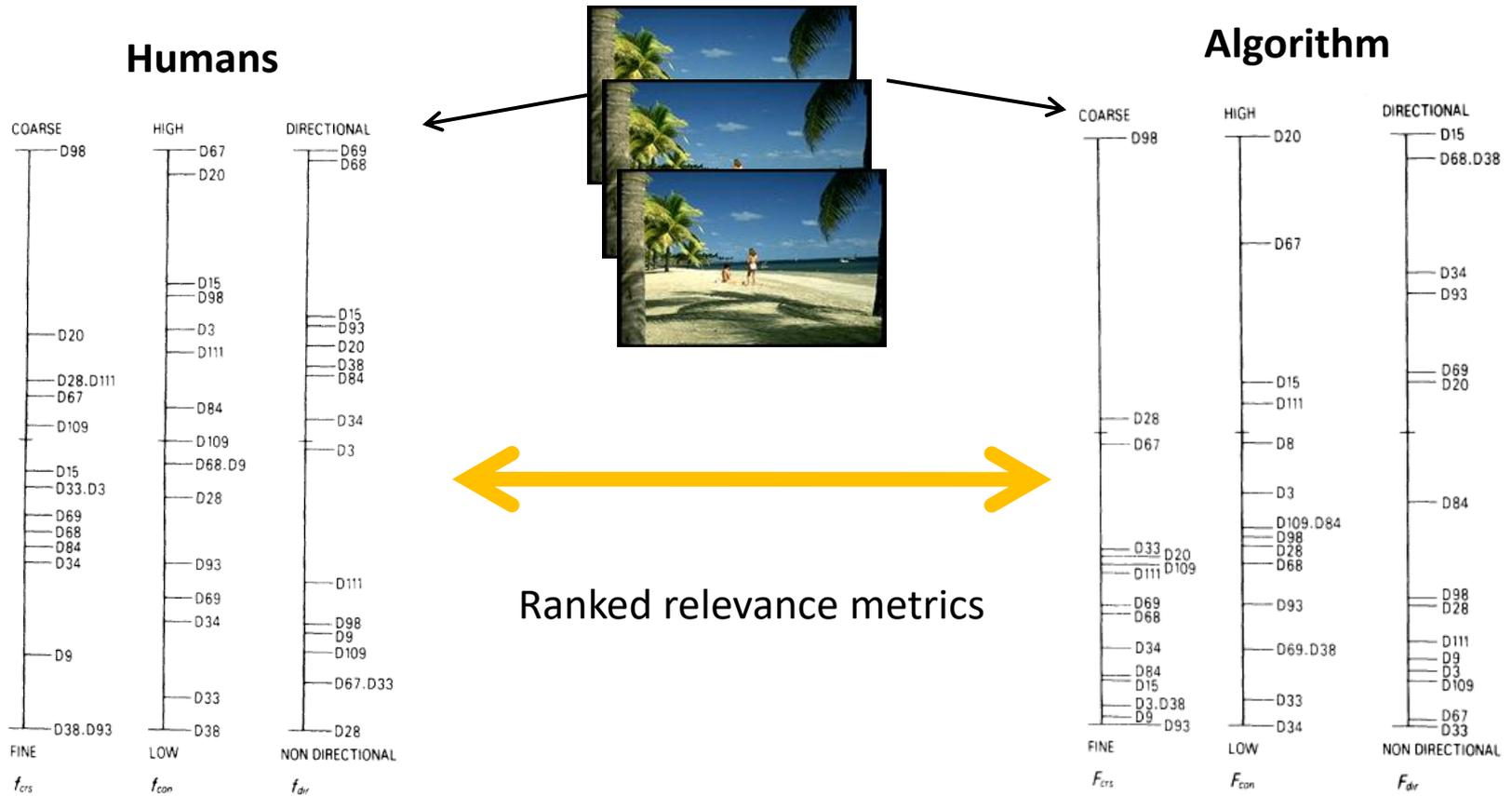
Psychological based textures (Tamura)

- **Coarseness** measures the size of the primitive elements forming the texture
- **Contrast** measures variation in gray levels between black and white
- **Directionality** measures the orientation of the texture
- **Line-likeness** measures the similarity of the texture to lines
- **Regularity** measures the repetitiveness of the texture pattern
- **Roughness** “we do not have any good ideas for describing the tactile sense of roughness”

Psychological based textures (Tamura)



Comparing psychological relevance to algorithms



Frequency based textures

- Frequency based texture decompose images according to their frequencies
 - Similar to audio filtering or color filter lenses
- The number of repetitions per area in a texture is related to the frequency of a texture
- Based on the Fourier Transform
- A set of 2 dimensional filters will decompose images into their natural frequencies

Edge detection



Edge detection

- Filter image with a low pass filter
- Apply vertical and horizontal filters to compute G_x and G_y :

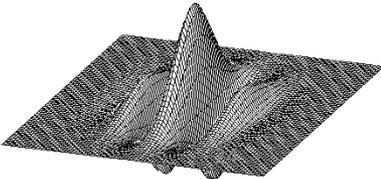
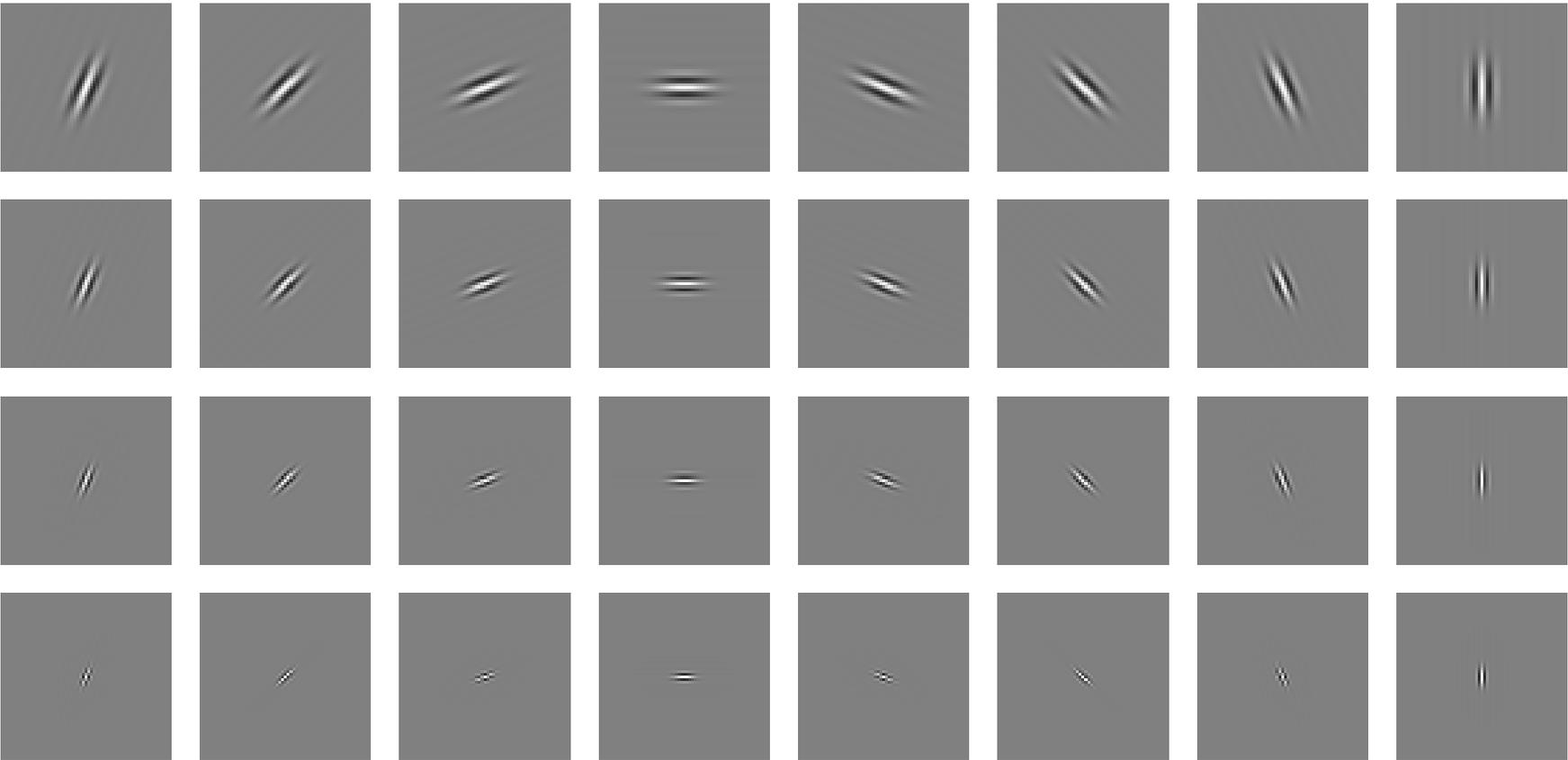
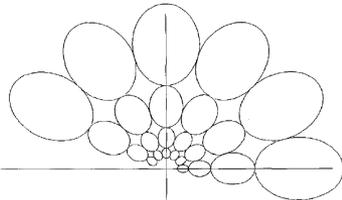
+1	+2	+1
0	0	0
-1	-2	-1

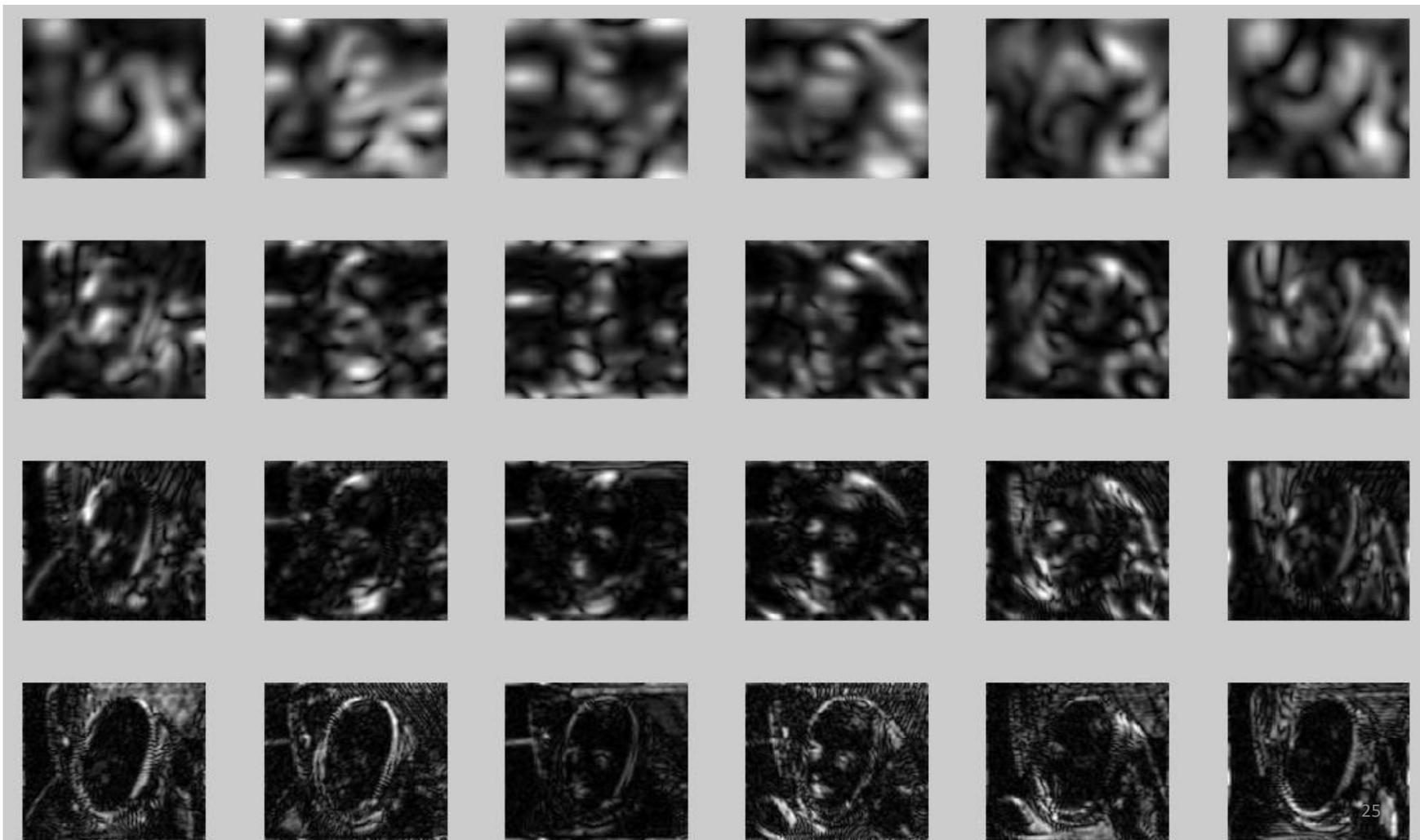
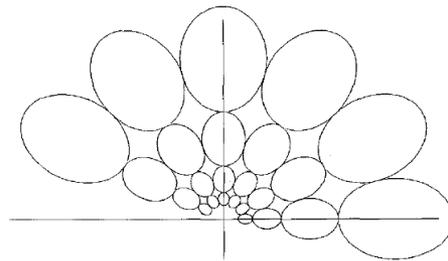
-1	0	+1
-2	0	+2
-1	0	+1

- Compute the gradients as $G = \sqrt{G_x^2 + G_y^2}$
 - Reduce it to one of the 4 possible directions (0°, 45°, 90°, 135°)

- Compute the orientation of the edges as: $\Theta = \arctan\left(\frac{G_y}{G_x}\right)$

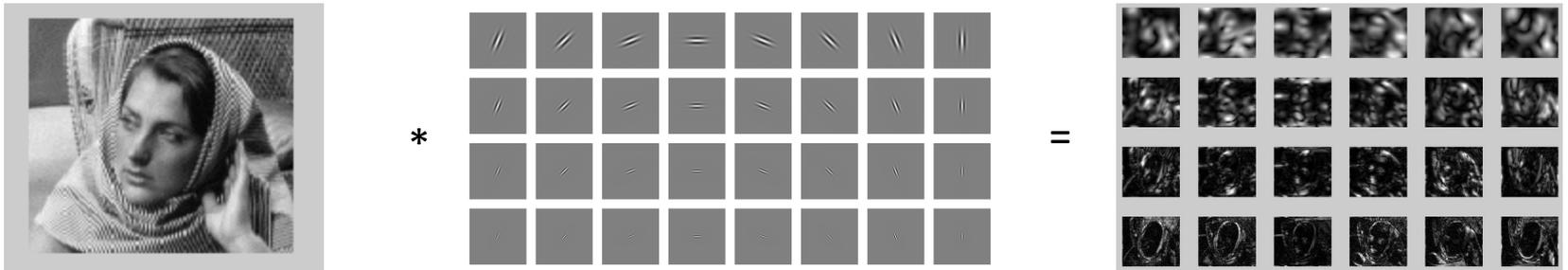
Gabor filters





Gabor texture feature

- Images are convolved (operator $*$) with each filter individually:



A widely used descriptor corresponds to the mean and variance of the output of each filter:

$$d_{texture} = (m_1, v_1, \dots, m_k, v_k)$$



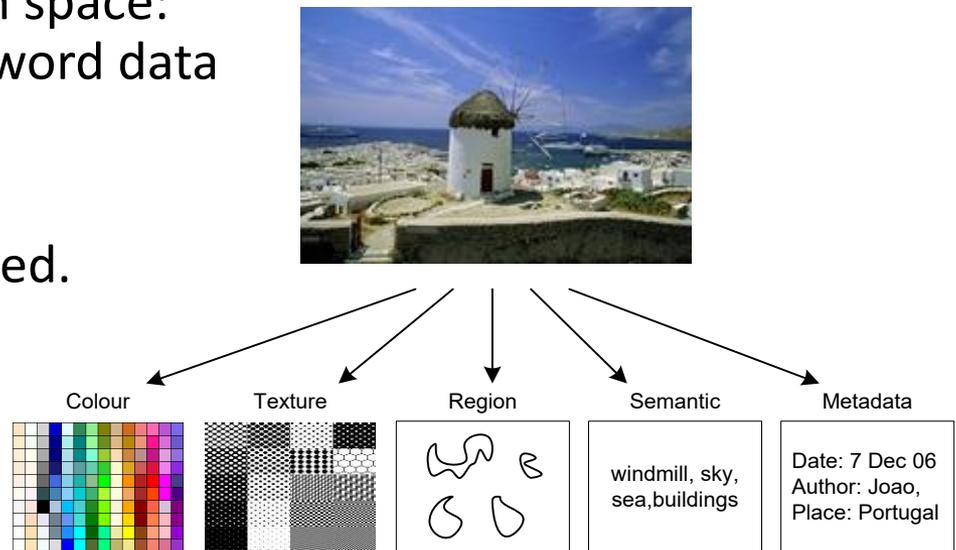
Multiple representations of the same data

- Documents are represented as the set of vectors

$$d = (d_{links}, d_{text}, d_{color}, d_{texture}, d_{metadata}, d_{tags}, \dots)$$

each one for a different search space:
text data, visual data, and keyword data
respectively.

- Other search spaces can be used.





Data representations

- Link data

$$d_{links} = (0,0, \dots, 0,1,0, \dots, 0,1,0, \dots, 0)$$

- High-dimensional data

- Sparse

- Bag of words

- Dense

- Color histograms and moments
 - Textures and edges

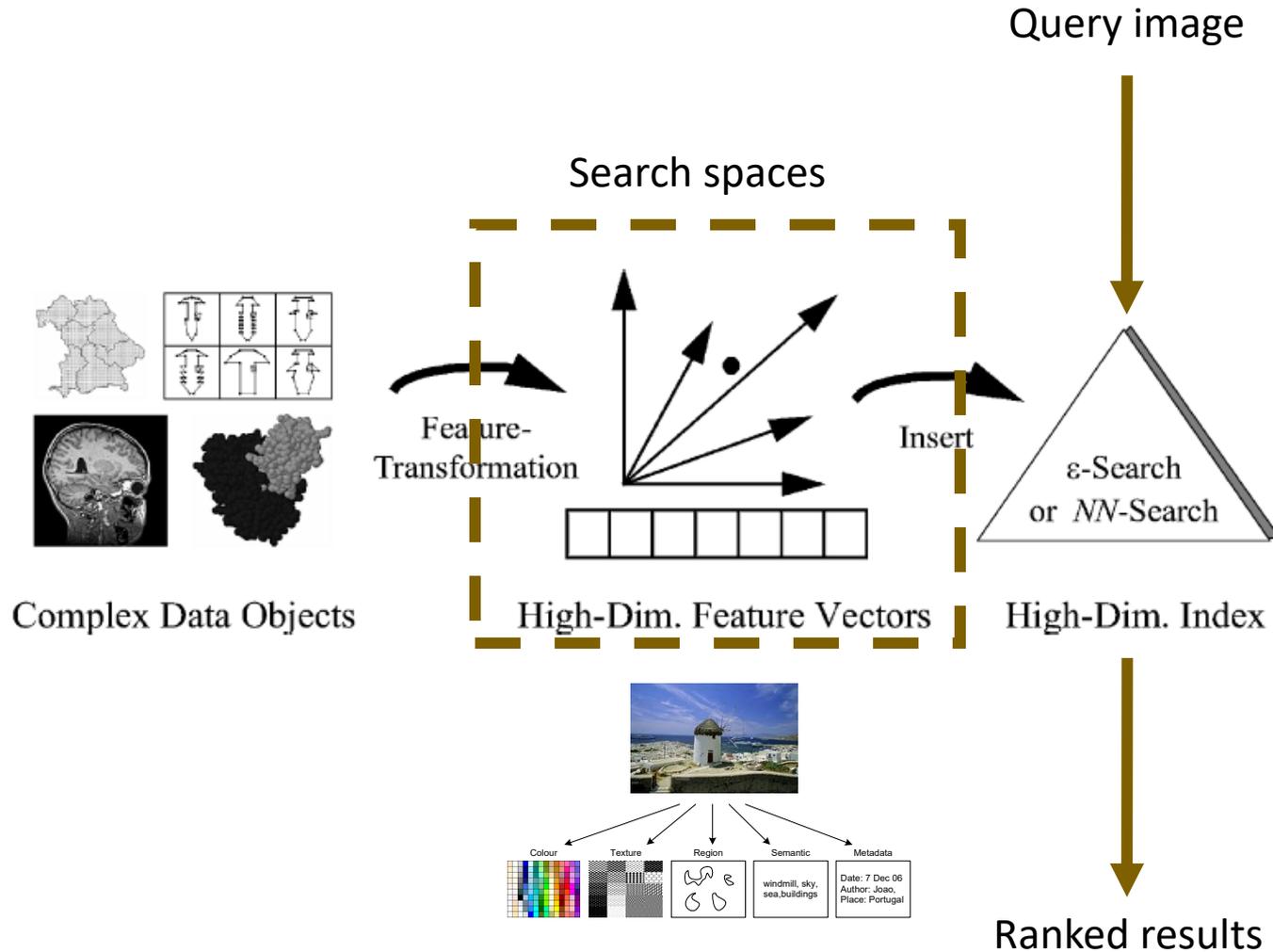
$$d_{bow} = (w_1, \dots, w_L, ng_1, \dots, ng_M)$$

$$d_{color} = (bin_1, bin_2, \dots, bin_k)$$

$$d_{texture} = (m_1, v_1, \dots, m_k, v_k)$$



Search high-dimensional spaces





Definition: metric spaces

- Let \mathcal{D} be an n dimensional space, where each data point is defined as

$$d \in \mathcal{D}: d = (d_1, \dots, d_n), \quad d_i \in \mathbb{R}$$

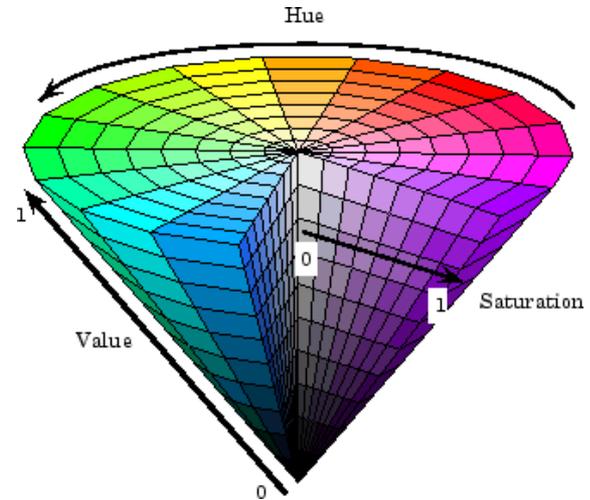
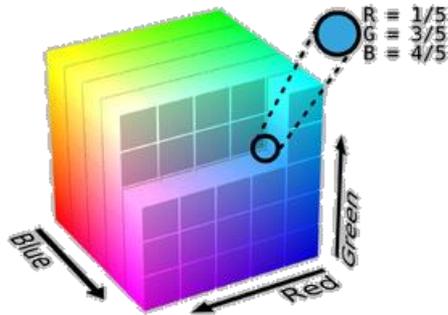
- The n dimensional space \mathcal{D} is a metric space iff exists a distance function $dist(a, b)$ in \mathcal{D} .
- A distance function has the following properties:
 - Non-negative: $dist(a, b) \geq 0 \quad \forall a, b \in \mathcal{D}$
 - Identity: if $dist(a, b) = 0$ then $a = b$
 - Symmetry: $dist(a, b) = dist(b, a) \quad \forall a, b \in \mathcal{D}$
 - Triangle inequality $dist(a, b) \leq dist(a, c) + dist(c, b) \quad \forall a, b, c \in \mathcal{D}$



Distance vs similarity

- Distances in a given search space must be meaningful.
- Distances are used as proxies for similarity.
 - $\text{distance} = 1 - \text{similarity}$
- Vector spaces and probability spaces are common spaces in Web search.
- The goal is that the similarity/distance between a query and candidate documents will reflect the relevance of the document to the user query.

Example: Distance in the RGB vs HSV color spaces

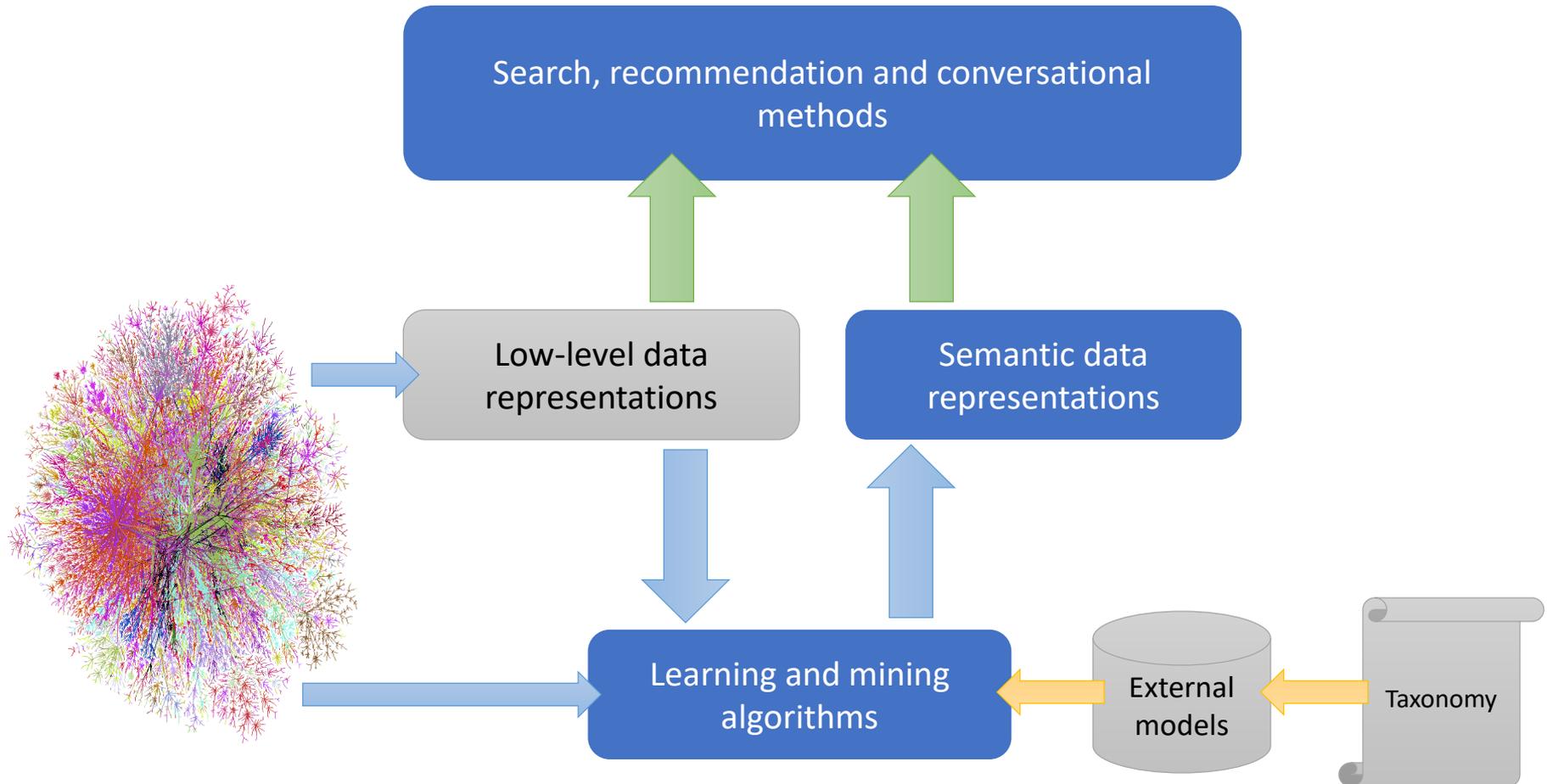


- Euclidean distance in the HSV color space is more meaningful!
 - Hue (H), the color type (such as red, green). It ranges from 0 to 360 degree.
 - Saturation (S) of the color ranges from 0 to 100%. Also sometimes it called the "purity".
 - Value (V), the Brightness (B) of the color ranges from 0 to 100%.

Searching Web content

- Processing real-world information is challenging!!!
- The aim is to search any unstructured data by its content
 - Textual, visual, audio, semantic, etc.
- Data contains very complex information patterns.
- Information needs can be very complex.
 - Queries can be *keywords, examples* or *questions*.
 - Finding related trends (consumption patterns)
 - Search images with text and vice-versa

Web Mining and Search course scope



- Multi-feature search
- Learning data representations
- Graph data (Labels and PageRank)
- Recommendation

Summary and readings

- Understanding the **diversity of Web data**
- Understanding the need to have different **data representations**
- **Spaces and similarity functions**
- References:
 - [Chapter 2](#): C. D. Manning, P. Raghavan and H. Schütze, “Introduction to Information Retrieval”, Cambridge University Press, 2008.
 - Hassaballah, M., Abdelmgeid, A. A., & Alshazly, H. A. (2016). [Image features detection, description and matching](#). In *Image Feature Detectors and Descriptors* (pp. 11-45). Springer, Cham.