

---

# Apprentissage Automatique et Fouille de données textuelles

Jean-Michel RENDERS

Xerox Research Center Europe (France)

AAFD'06

# Plan Global

---

- Introduction :
  - Fouille de textes
  - Spécificité des données textuelles
- Approche numéro 1 : méthodes à noyaux
  - Philosophie des méthodes à noyaux
  - Noyaux pour les données textuelles
- Approche numéro 2 : modèles génératifs
  - Génératif versus discriminatif – semi-supervisé
  - Modèles graphiques à variables latentes
  - Exemples : NB, PLSA, LDA, HPLSA
- Perspectives « récentes »

# Fouille de Textes?

---

- Sens strict : très rare
  - Sens large: contient une panoplie de sous-tâches
    - Recherche d'information (IR->QA)
    - Analyse sémantique
    - Catégorisation, Clustering
    - Extraction d'information → population d'ontologie
    - Focalisation utilisateur: navigation, visualisation, résumé adapté, traduction, ...
- Souvent précédée de tâches de pré-traitement linguistique (jusqu'à l'analyse syntaxique et le tagging)  
... elles-mêmes appelées Fouille de textes!

# Spécificités du Texte

---

- Qu'est-ce qu'une observation?
  - Objet d'étude à différents niveaux de granularité (mot, phrase, section, document, corpus, mais aussi utilisateur, communauté)
- Lien entre forme et fond
  - Paradoxe structuré – non structuré
  - Importance d'un background knowledge
  - Redondance (cfr. Synonymie) et ambiguïté (cfr. Polysémie)

# Cas particulier

---

## ■ Cas d'école le plus fréquent

- Objet d'étude: document
- Attributs: mots

## ■ Propriétés:

- Attributs: polysémie, synonymie, structuration hiérarchique, dépendance ordonnée, attributs composés
- Documents: polythématicité, structuration des classes, appartenance floue

# Polythématicité



“This website documents a period in the life of Arundhati Roy when she became a world-renowned writer, first as the author of the Booker prize-winning *The God of Small Things*, and then as a campaigning activist and inspiration to all those who seek to support the poor and oppressed and to speak up against the power that be.”

# Approach 1 – Kernel Methods

---

- What's the philosophy of Kernel Methods?
- How to use Kernel Methods in Learning tasks?
- Kernels for text (BOW, latent concept, string, word sequence, tree and Fisher Kernels)
- Applications to NLP tasks

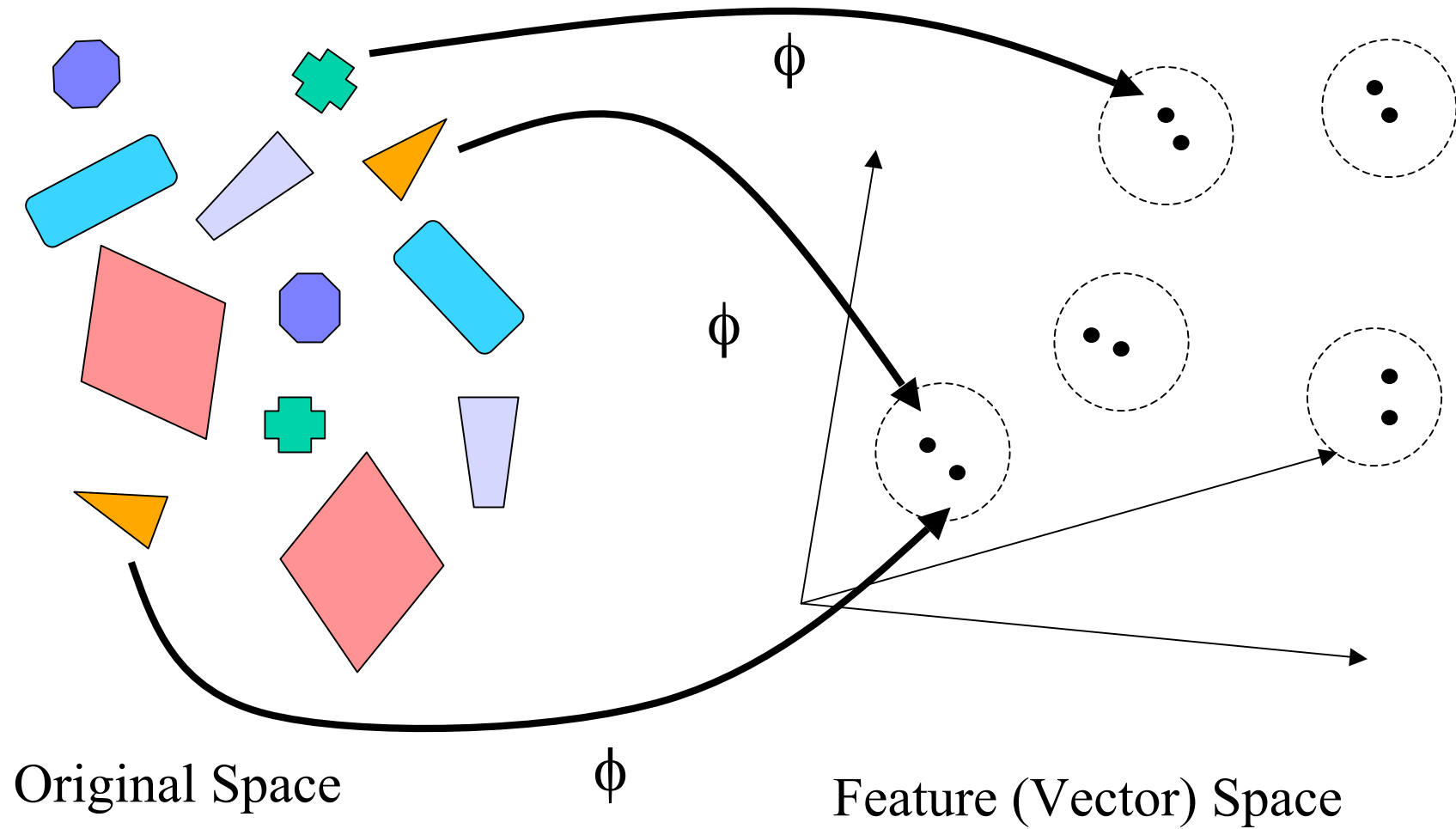
# Kernel Methods : intuitive idea

---

- Find a mapping  $\phi$  such that, in the new space, problem solving is easier (e.g. linear)
- The *kernel* represents the similarity between two objects (documents, terms, ...), defined as the dot-product in this new vector space
- But the mapping is left implicit
- Easy generalization of a lot of dot-product (or distance) based pattern recognition algorithms



# Kernel Methods : the mapping



# Kernel : more formal definition

---

- A kernel  $k(x,y)$ 
  - is a similarity measure
  - defined by an implicit mapping  $\phi$ ,
  - from the original space to a vector space (feature space)
  - such that:  $k(x,y)=\phi(x)\cdot\phi(y)$
- This similarity measure and the mapping include:
  - Invariance or other a priori knowledge
  - Simpler structure (linear representation of the data)
  - The class of functions the solution is taken from
  - Possibly infinite dimension (hypothesis space for learning)
  - ... but still computational efficiency when computing  $k(x,y)$

# Benefits from kernels

---

- Generalizes (nonlinearly) pattern recognition algorithms in clustering, classification, density estimation, ...
  - When these algorithms are dot-product based, by replacing the dot product  $(x \cdot y)$  by  $k(x,y) = \phi(x) \cdot \phi(y)$   
e.g.: linear discriminant analysis, logistic regression, perceptron, SOM, PCA, ICA, ...  
NM. This often implies to work with the “dual” form of the algo.
  - When these algorithms are distance-based, by replacing  $d(x,y)$  by  $k(x,x) + k(y,y) - 2k(x,y)$
- Freedom of choosing  $\phi$  implies a large variety of learning algorithms

# Valid Kernels

---

- The function  $k(x,y)$  is a valid kernel, if there exists a mapping  $\phi$  into a vector space (with a dot-product) such that  $k$  can be expressed as  $k(x,y)=\phi(x)\cdot\phi(y)$
- Theorem:  $k(x,y)$  is a valid kernel if  $k$  is positive definite and symmetric (Mercer Kernel)
  - A function is P.D. if  $\int K(\mathbf{x},\mathbf{y})f(\mathbf{x})f(\mathbf{y})d\mathbf{x}d\mathbf{y} \geq 0 \quad \forall f \in L_2$
  - In other words, the Gram matrix  $\mathbf{K}$  (whose elements are  $k(x_i,x_j)$ ) must be positive definite for all  $x_i, x_j$  of the input space
  - One possible choice of  $\phi(x)$ :  $k(\cdot,x)$  (maps a point  $x$  to a function  $k(\cdot,x) \rightarrow$  feature space with infinite dimension!)

# Example of Kernels (I)

---

- Polynomial Kernels:  $k(x,y)=(x\cdot y)^d$ 
  - Assume we know most information is contained in monomials (e.g. multiword terms) of degree  $d$  (e.g.  $d=2$ :  $x_1^2, x_2^2, x_1x_2$ )
  - Theorem: the (implicit) feature space contains all possible monomials of degree  $d$  (ex:  $n=250$ ;  $d=5$ ;  $\dim F=10^{10}$ )
  - But kernel computation is only marginally more complex than standard dot product!
  - For  $k(x,y)=(x\cdot y+1)^d$ , the (implicit) feature space contains all possible monomials up to degree  $d$  !

# The Kernel Gram Matrix

---

- With KM-based learning, the **sole** information used from the training data set is the Kernel Gram Matrix

$$K_{training} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_m) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \dots & k(\mathbf{x}_2, \mathbf{x}_m) \\ \dots & \dots & \dots & \dots \\ k(\mathbf{x}_m, \mathbf{x}_1) & k(\mathbf{x}_m, \mathbf{x}_2) & \dots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

- If the kernel is valid,  $K$  is symmetric definite-positive .

# How to build new kernels

---

- Kernel combinations, preserving *validity*:

$$K(\mathbf{x}, \mathbf{y}) = \lambda K_1(\mathbf{x}, \mathbf{y}) + (1 - \lambda) K_2(\mathbf{x}, \mathbf{y}) \quad 0 \leq \lambda \leq 1$$

$$K(\mathbf{x}, \mathbf{y}) = a \cdot K_1(\mathbf{x}, \mathbf{y}) \quad a > 0$$

$$K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) \cdot K_2(\mathbf{x}, \mathbf{y})$$

$$K(\mathbf{x}, \mathbf{y}) = f(x) \cdot f(y) \quad f \text{ is real-valued function}$$

$$K(\mathbf{x}, \mathbf{y}) = K_3(\ddot{\mathbf{o}}(\mathbf{x}), \ddot{\mathbf{o}}(\mathbf{y}))$$

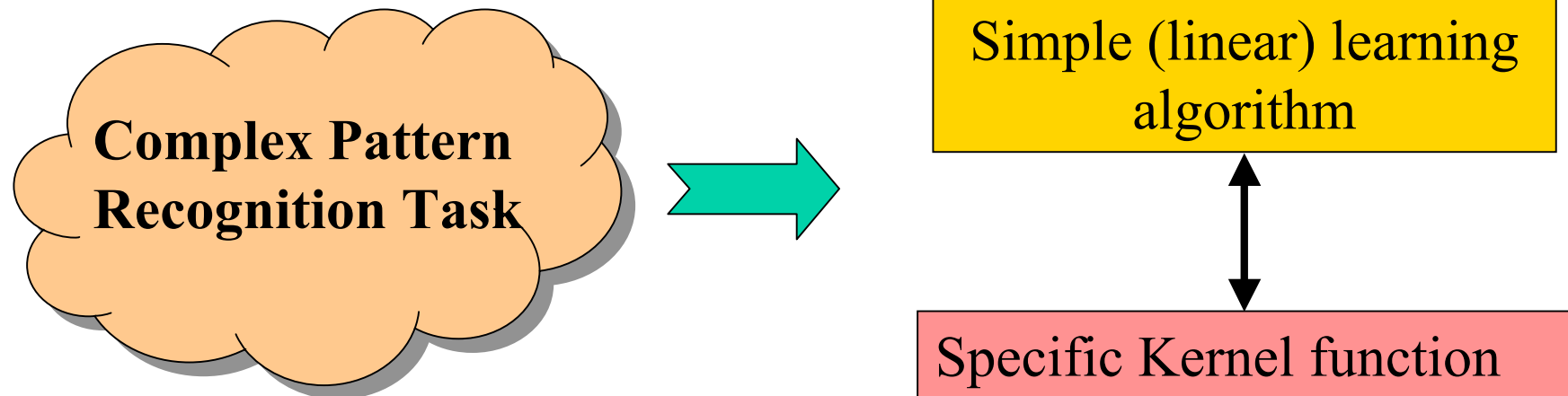
$$K(\mathbf{x}, \mathbf{y}) = \mathbf{x}' P \mathbf{y} \quad P \text{ symmetric definite positive}$$

$$K(\mathbf{x}, \mathbf{y}) = \frac{K_1(\mathbf{x}, \mathbf{y})}{\sqrt{K_1(\mathbf{x}, \mathbf{x})} \sqrt{K_1(\mathbf{y}, \mathbf{y})}}$$

# Kernels and Learning

---

- In Kernel-based learning algorithms, problem solving is now decoupled into:
  - A general purpose learning algorithm (e.g. SVM, PCA, ...) – Often linear algorithm (well-funded, robustness, ...)
  - A problem specific kernel





# Learning in the feature space: Issues

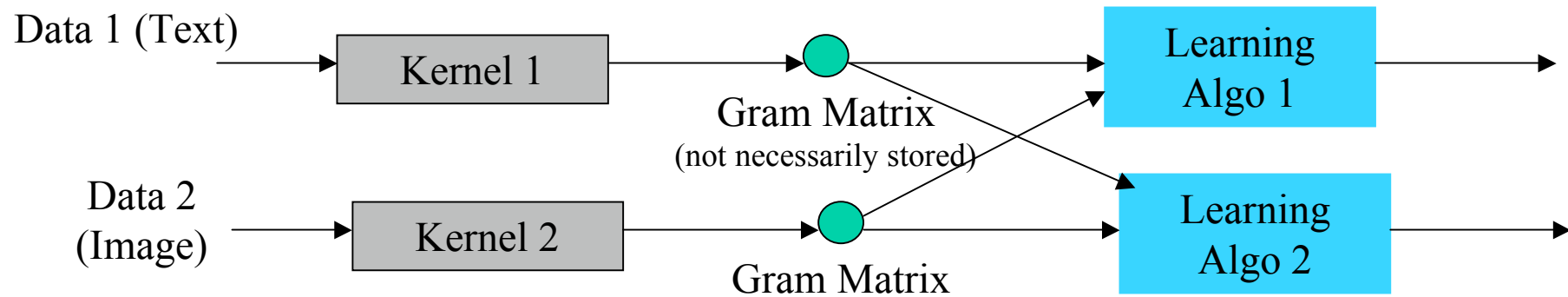
---

- High dimensionality allows to render flat complex patterns by “explosion”
  - Computational issue, solved by designing kernels (efficiency in space and time)
  - Statistical issue (generalization), solved by the learning algorithm and also by the kernel
    - e.g. SVM, solving this complexity problem by maximizing the margin and the dual formulation
- E.g. RBF-kernel, playing with the  $\sigma$  parameter
- With adequate learning algorithms and kernels, high dimensionality is no longer an issue

# Current Synthesis

---

- Modularity and re-usability
  - Same kernel ,different learning algorithms
  - Different kernels, same learning algorithms
- This presentation is allowed to focus only on designing kernels for textual data



# Agenda

---

- What's the philosophy of Kernel Methods?
- How to use Kernels Methods in Learning tasks?
- Kernels for text (BOW, latent concept, string, word sequence, tree and Fisher Kernels)
- Applications to NLP tasks

# Kernels for texts

---

## ■ Similarity between documents?

- Seen as 'bag of words' : dot product or polynomial kernels (multi-words)
- Seen as set of concepts : GVSM kernels, Kernel LSI (or Kernel PCA), Kernel ICA, ...possibly multilingual
- Seen as string of characters: string kernels
- Seen as string of terms/concepts: word sequence kernels
- Seen as trees (dependency or parsing trees): tree kernels
- Seen as the realization of probability distribution (generative model)

# Strategies of Design

---

- Kernel as a way to encode prior information
  - Invariance: synonymy, document length, ...
  - Linguistic processing: word normalisation, semantics, stopwords, weighting scheme, ...
- Convolution Kernels: text is a recursively-defined data structure. How to build “global” kernels from local (atomic level) kernels?
- Generative model-based kernels: the “topology” of the problem will be translated into a kernel function (cfr. Mahalanobis)

# Strategies of Design

---

- Kernel as a way to encode prior information
  - Invariance: synonymy, document length, ...
  - Linguistic processing: word normalisation, semantics, stopwords, weighting scheme, ...
- Convolution Kernels: text is a recursively-defined data structure. How to build “global” kernels from local (atomic level) kernels?
- Generative model-based kernels: the “topology” of the problem will be translated into a kernel function

# 'Bag of words' kernels (I)

---

- Document seen as a vector  $\mathbf{d}$ , indexed by all the elements of a (controlled) dictionary. The entry is equal to the number of occurrences.
- A training corpus is therefore represented by a Term-Document matrix,  
noted  $\mathbf{D}=[\mathbf{d}_1 \ \mathbf{d}_2 \ \dots \ \mathbf{d}_{m-1} \ \mathbf{d}_m]$
- The “nature” of word: will be discussed later
- From this basic representation, we will apply a sequence of successive embeddings, resulting in a global (valid) kernel with all desired properties

# BOW kernels (II)

---

## ■ Properties:

- All order information is lost (syntactical relationships, local context, ...)
- Feature space has dimension N (size of the dictionary)

## ■ Similarity is basically defined by:

$$k(d_1, d_2) = \mathbf{d}_1 \cdot \mathbf{d}_2 = \mathbf{d}_1^t \cdot \mathbf{d}_2$$

or, normalized (cosine similarity):

$$\hat{k}(d_1, d_2) = \frac{k(d_1, d_2)}{\sqrt{k(d_1, d_1) \cdot k(d_2, d_2)}}$$

- Efficiency provided by sparsity (and sparse dot-product algo):  $O(|d_1| + |d_2|)$



# 'Bag of words' kernels: enhancements

---

- The choice of indexing terms:
  - Exploit linguistic enhancements:
    - Lemma / Morpheme & stem
    - Disambiguated lemma (lemma+POS)
    - Noun Phrase (or useful collocation, n-grams)
    - Named entity (with type)
  - Exploit IR lessons
    - Stopword removal
    - Feature selection based on frequency
    - Weighting schemes (e.g. idf )
    - Semantic enrichment by term-term similarity matrix  $\mathbf{Q}$  (positive definite):  
 $k(\mathbf{d}_1, \mathbf{d}_2) = \phi(\mathbf{d}_1)^t \cdot \mathbf{Q} \cdot \phi(\mathbf{d}_2)$
- NB. Using polynomial kernels up to degree  $p$ , is a natural and efficient way of considering all (up-to-)  $p$ -grams (with different weights actually), but order is not taken into account (“sinking ships” is the same as “shipping sinks”)

# Semantic Smoothing Kernels

---

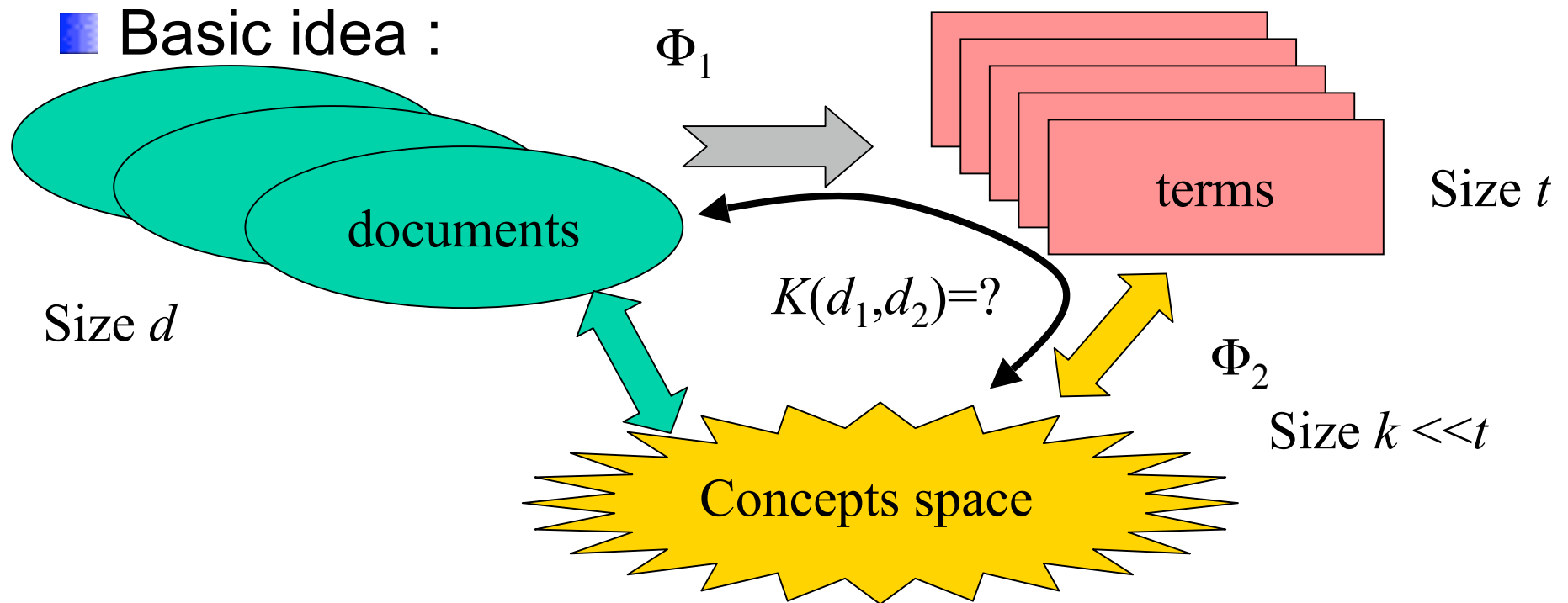
- Synonymy and other term relationships:
  - GVSM Kernel: the term-term co-occurrence matrix ( $\mathbf{DD}^t$ ) is used in the kernel:  $k(d_1, d_2) = \mathbf{d}_1^t \cdot (\mathbf{D} \cdot \mathbf{D}^t) \cdot \mathbf{d}_2$
  - The completely kernelized version of GVSM is:
    - The training kernel matrix  $\mathbf{K} (= \mathbf{D}^t \cdot \mathbf{D}) \rightarrow \mathbf{K}^2 (m \times m)$
    - The kernel vector of a new document  $d$  vs the training documents :  $\mathbf{t} \rightarrow \mathbf{K} \cdot \mathbf{t} (m \times 1)$
    - The initial  $\mathbf{K}$  could be a polynomial kernel (GVSM on multi-words terms)
  - Variants: One can use
    - a shorter context than the document to compute term-term similarity (term-context matrix)
    - Another measure than the number of co-occurrences to compute the similarity (e.g. Mutual information, ...)
  - Can be generalised to  $\mathbf{K}^n$  (or a weighted combination of  $\mathbf{K}^1 \mathbf{K}^2 \dots \mathbf{K}^n$  cfr. Diffusion kernels later), but is  $\mathbf{K}^n$  less and less sparse! Interpretation as sum over paths of length  $2n$ .

# Semantic Smoothing Kernels

---

- Can use other term-term similarity matrix than  $\mathbf{DD}^t$ ; e.g. a similarity matrix derived from the Wordnet thesaurus, where the similarity between two terms is defined as:
  - the inverse of the length of the path connecting the two terms in the hierarchical hyper/hyponymy tree.
  - A similarity measure for nodes on a tree (feature space indexed by each node  $n$  of the tree, with  $\phi_n(\text{term } x)$  if term  $x$  is the class represented by  $n$  or “under”  $n$ ), so that the similarity is the number of common ancestors (including the node of the class itself).
- With semantic smoothing, 2 documents can be similar even if they don't share common words.

# Latent concept Kernels



# Latent concept Kernels

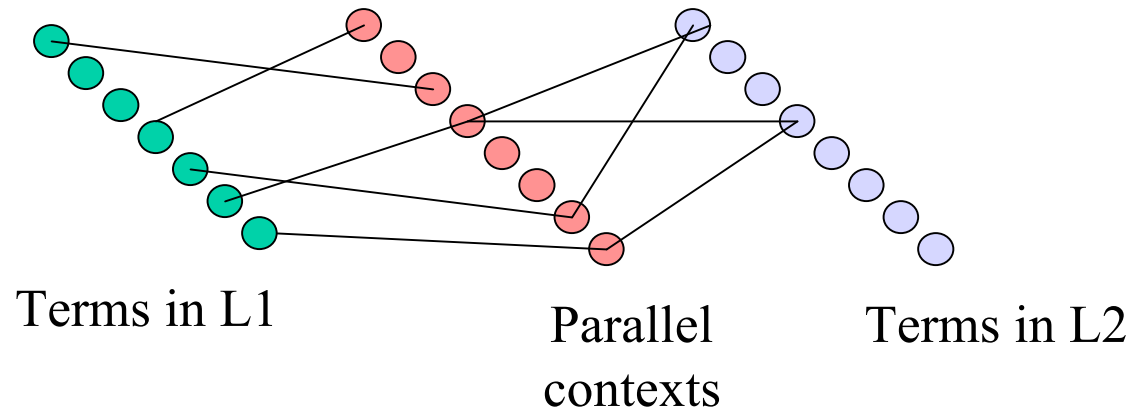
---

- $k(d_1, d_2) = \phi(\mathbf{d}_1)^t \cdot \mathbf{P}^t \cdot \mathbf{P} \cdot \phi(\mathbf{d}_2)$ ,
  - where  $\mathbf{P}$  is a (linear) projection operator
    - From Term Space
    - to Concept Space
- Working with (latent) concepts provides:
  - Robustness to polysemy, synonymy, style, ...
  - Cross-lingual bridge
  - Natural Dimension Reduction
- But, how to choose  $\mathbf{P}$  and how to define (extract) the latent concept space? Ex: Use PCA : the concepts are nothing else than the principal components.

# Why multilingualism helps ...

---

## ■ Graphically:



- Concatenating both representations will force language-independent concept: each language imposes constraints on the other
- Searching for maximally correlated projections of paired observations (CCA) has a sense, semantically speaking

# Diffusion Kernels

---

- Recursive dual definition of the semantic smoothing:

$$\mathbf{K} = \mathbf{D}'(\mathbf{I} + u\mathbf{Q})\mathbf{D}$$

$$\mathbf{Q} = \mathbf{D}(\mathbf{I} + v\mathbf{K})\mathbf{D}'$$

NB.  $u=v=0 \rightarrow$  standard BOW;  $v=0 \rightarrow$  GVSM

- Let  $\mathbf{B} = \mathbf{D}'\mathbf{D}$  (standard BOW kernel);  $\mathbf{G} = \mathbf{D}\mathbf{D}'$

- If  $u=v$ , The solution is the “Von Neumann diffusion kernel”

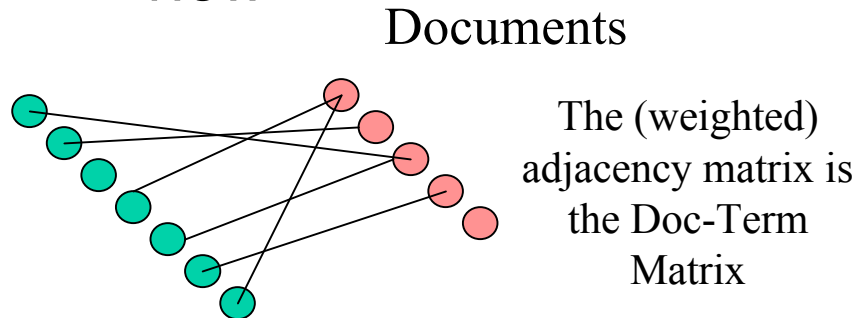
- $\mathbf{K} = \mathbf{B} \cdot (\mathbf{I} + u\mathbf{B} + u^2\mathbf{B}^2 + \dots) = \mathbf{B}(\mathbf{I} - u\mathbf{B})^{-1}$  and  $\mathbf{Q} = \mathbf{G}(\mathbf{I} - u\mathbf{G})^{-1}$  [only if  $u < \|\mathbf{B}\|^{-1}$ ]

- Can be extended, with a faster decay, to exponential diffusion kernel:

- $\mathbf{K} = \mathbf{B} \cdot \exp(u\mathbf{B})$  and  $\mathbf{Q} = \exp(u\mathbf{G})$

# Graphical Interpretation

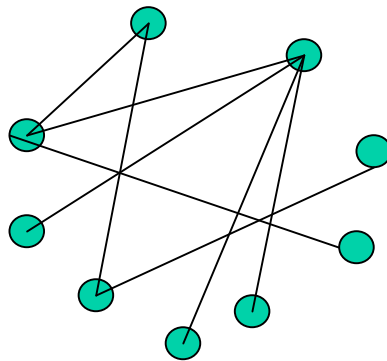
- These diffusion kernels correspond to defining similarities between nodes in a graph, specifying only the myopic view



Terms

Or

By aggregation, the (weighted) adjacency matrix is the term-term similarity matrix  $\mathbf{G}$



*Diffusion kernels corresponds to considering all paths of length 1, 2, 3, 4 ... linking 2 nodes and summing the product of local similarities, with different decay strategies*

*It is in some way similar to KPCA by just “rescaling” the eigenvalues of the basic Kernel matrix (decreasing the lowest ones)*



# Strategies of Design

---

- Kernel as a way to encode prior information
  - Invariance: synonymy, document length, ...
  - Linguistic processing: word normalisation, semantics, stopwords, weighting scheme, ...

■ **Convolution Kernels:** text is a recursively-defined data structure. How to build “global” kernels from local (atomic level) kernels?

- Generative model-based kernels: the “topology” of the problem will be translated into a kernel function

# Sequence kernels

---

## ■ Consider a document as:

- A sequence of characters (string)
- A sequence of tokens (or stems or lemmas)
- A paired sequence (POS+lemma)
- A sequence of concepts

- A tree (parsing tree)
  - A dependency graph
- (later)

## ■ Sequence kernels → order has importance

- Kernels on string/sequence : counting the subsequences two objects have in common ... but various ways of counting
- Contiguity is necessary (p-spectrum kernels)
- Contiguity is not necessary (subsequence kernels)
- Contiguity is penalised (gap-weighted subsequence kernels)

# String and Sequence

---

- Just a matter of convention:
  - String matching: implies contiguity
  - Sequence matching : only implies order

# Gap-weighted subsequence kernels

---

- Feature space indexed by all elements of  $\Sigma^p$
- $\phi_u(s)$  = sum of weights of occurrences of the p-gram  $u$  as a (non-contiguous) subsequence of  $s$ , the weight being length penalizing:  $\lambda^{\text{length}(u)}$  [NB: length includes both matching symbols and gaps]
- Example:
  - D1 : ATCGTAGACTGTC
  - D2 : GACTATGC
  - $(D1)_{\text{CAT}} = 2\lambda^8 + 2\lambda^{10}$  and  $(D2)_{\text{CAT}} = \lambda^4$
  - $k(D1, D2)_{\text{CAT}} = 2\lambda^{12} + 2\lambda^{14}$
- Naturally built as a dot product  $\rightarrow$  valid kernel
- For alphabet of size 80, there are 512000 trigrams
- For alphabet of size 26, there are  $12 \cdot 10^6$  5-grams

# Gap-weighted subsequence kernels

---

- Hard to perform explicit expansion and dot-product!
- Efficient recursive formulation (dynamic programming –like), whose complexity is  $O(k \cdot |D1| \cdot |D2|)$
- Normalization (doc length independence)

$$\hat{k}(d_1, d_2) = \frac{k(d_1, d_2)}{\sqrt{k(d_1, d_1) \cdot k(d_2, d_2)}}$$

# Word Sequence Kernels (I)

---

- Here “words” are considered as symbols
  - Meaningful symbols → more relevant matching
  - Linguistic preprocessing can be applied to improve performance
  - Shorter sequence sizes → improved computation time
  - But increased sparsity (documents are more : “orthogonal”)
  - Intermediate step: syllable kernel (indirectly realizes some low-level stemming and morphological decomposition)
- Motivation : the noisy stemming hypothesis (important N-grams approximate stems), confirmed experimentally in a categorization task

4-GRAMS

P	R	I	C	E	O	F	C	O	R	N	E	X	P	O	R	T	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

5-GRAMS

I	S	H	A	J	O	I	N	T	V	E	N	T	U	R	E	W
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

# Word Sequence Kernels (II)

---

- Link between Word Sequence Kernels and other methods:
  - For  $k=1$ , WSK is equivalent to basic “Bag Of Words” approach
  - For  $\lambda=1$ , close relation to polynomial kernel of degree  $k$ , WSK takes order into account
- Extension of WSK:
  - Symbol dependant decay factors (way to introduce IDF concept, dependence on the POS, stop words)
  - Different decay factors for gaps and matches (e.g.  $\lambda_{\text{noun}} < \lambda_{\text{adj}}$  when gap;  $\lambda_{\text{noun}} > \lambda_{\text{adj}}$  when match)
  - Soft matching of symbols (e.g. based on thesaurus, or on dictionary if we want cross-lingual kernels)

# Trie-based kernels

---

- An alternative to DP based on string matching techniques
- TRIE= Retrieval Tree (cfr. Prefix tree) = tree whose internal nodes have their children indexed by  $\Sigma$ .
- Suppose  $F = \Sigma^p$  : the leaves of a complete  $p$ -trie are the indices of the feature space
- Basic algorithm:
  - Generate all substrings  $s(i:j)$  satisfying initial criteria; idem for  $t$ .
  - Distribute the  $s$ -associated list down from root to leaf (depth-first)
  - Distribute the  $t$ -associated list down from root to leaf taking into account the distribution of  $s$ -list (pruning)
  - Compute the product at the leaves and sum over the leaves
- Key points: in steps (2) and (3), not all the leaves will be populated (else complexity would be  $O(|\Sigma^p|)$  ... you need not build the trie explicitly!



# Tree Kernels

---

- Application: categorization [one doc=one tree], parsing (desambiguation) [one doc = multiple trees]
- Tree kernels constitute a particular case of more general kernels defined on discrete structure (convolution kernels). Intuitively, the philosophy is
  - to split the structured objects in parts,
  - to define a kernel on the “atoms” and a way to recursively combine kernel over parts to get the kernel over the whole.

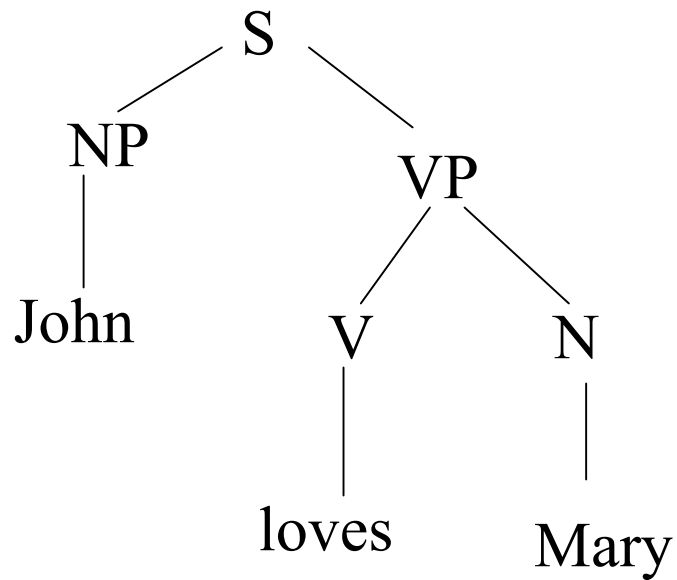
# Fundamentals of Tree kernels

---

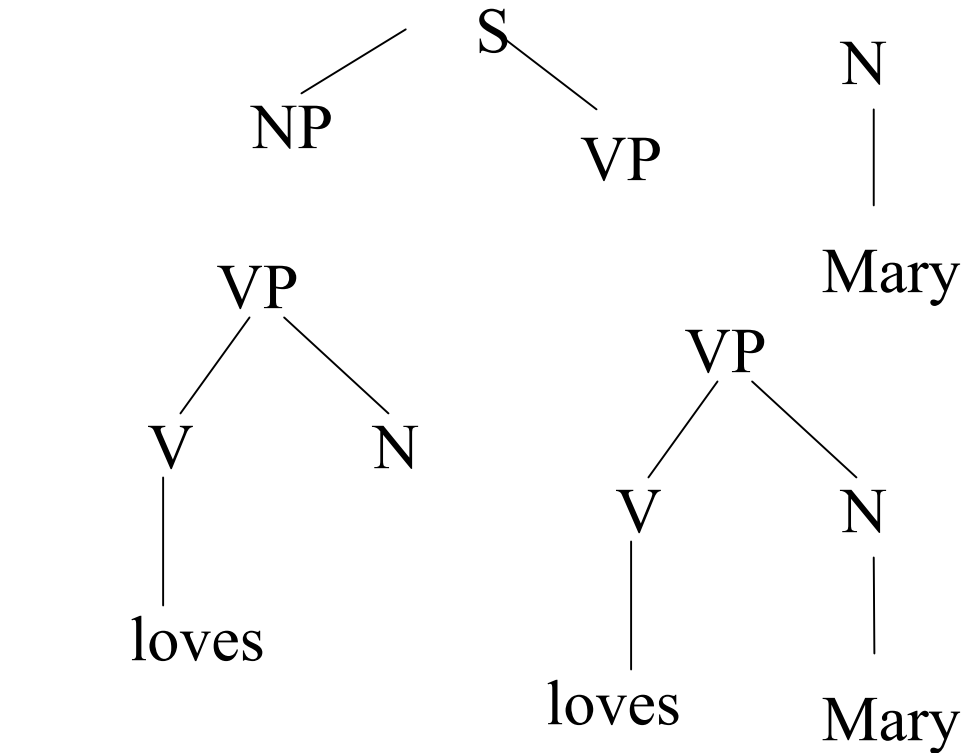
- Feature space definition: one feature for each possible proper subtree in the training data; feature value = number of occurrences
- A subtree is defined as any part of the tree which includes more than one node, with the restriction there is no “partial” rule production allowed.

# Tree Kernels : example

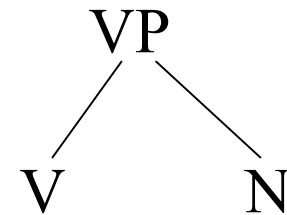
## ■ Example :



A Parse Tree



... a few among the  
many subtrees of  
this tree!



# Tree Kernels : algorithm

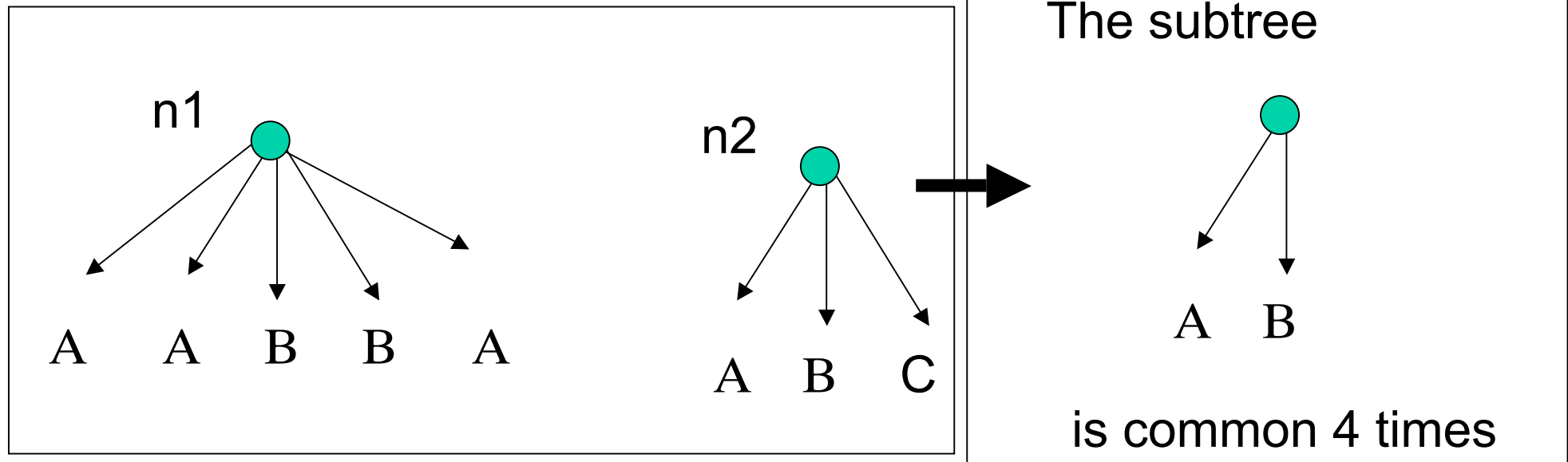
---

- Kernel = dot product in this high dimensional feature space
- Once again, there is an efficient recursive algorithm (in polynomial time, not exponential!)
- Basically, it compares the production of all possible pairs of nodes  $(n_1, n_2)$  ( $n_1 \in T_1, n_2 \in T_2$ ); if the production is the same, the number of common subtrees rooted at both  $n_1$  and  $n_2$  is computed recursively, considering the number of common subtrees rooted at the common children
- Formally, let  $k_{\text{co-rooted}}(n_1, n_2)$  = number of common subtrees rooted at both  $n_1$  and  $n_2$

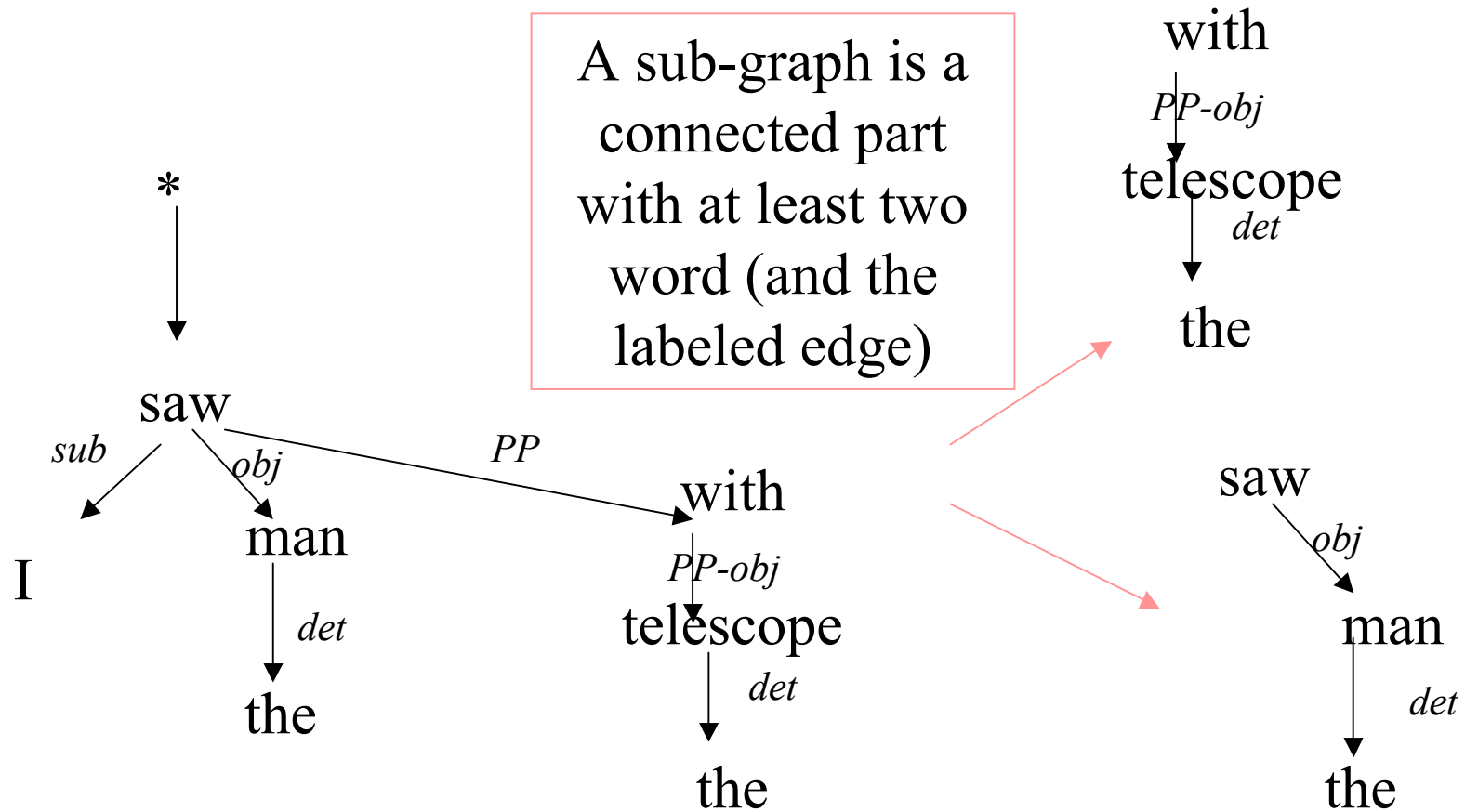
$$k(T_1, T_2) = \sum_{n_1 \in T_1} \sum_{n_2 \in T_2} k_{\text{co-rooted}}(n_1, n_2)$$

# Variant for labeled ordered tree

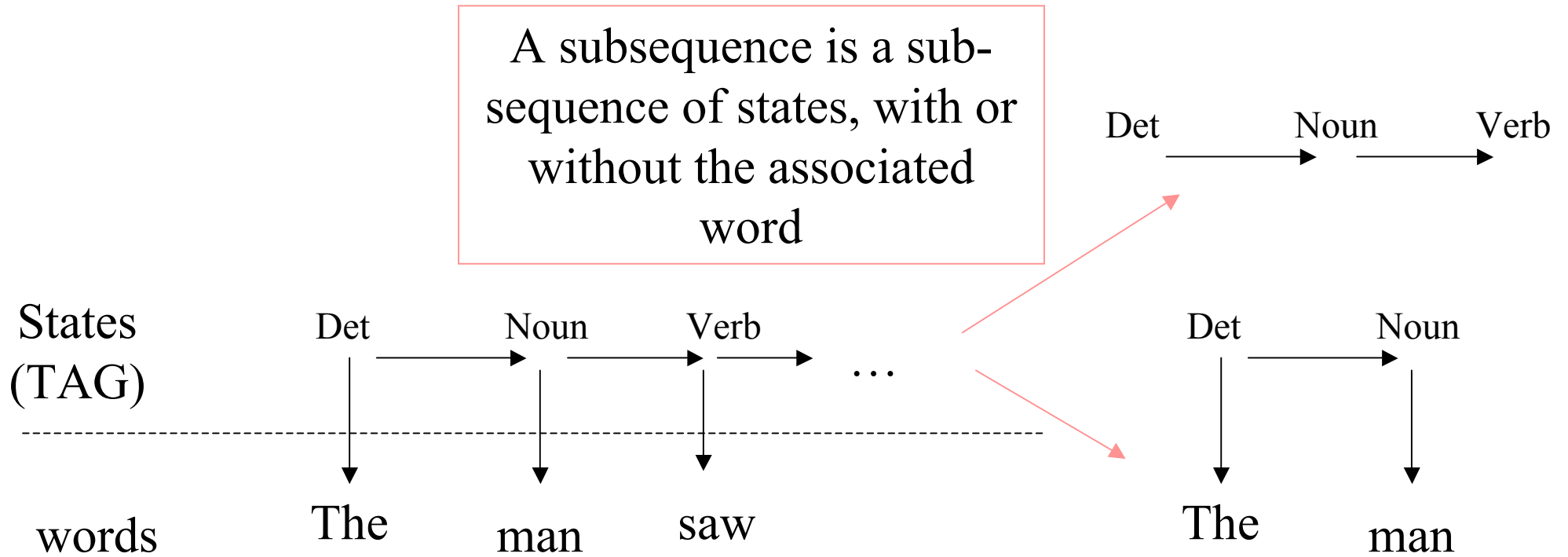
- Example: dealing with html/xml documents
- Extension to deal with:
  - Partially equal production
  - Children with same labels
  - ... but order is important



# Dependency Graph Kernel



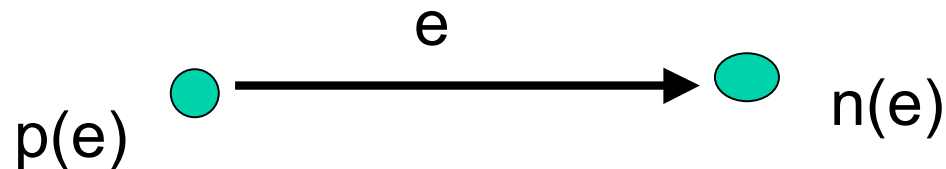
# Paired sequence kernel



# Graph kernels based on Common Walks

---

- Walk = (possibly infinite) sequence of labels obtained by following edges on the graph
- Path = walk with no vertex visited twice
- Important concept: direct product of two graphs  $G1 \times G2$ 
  - $V(G1 \times G2) = \{(v1, v2), v1 \text{ and } v2: \text{ same labels}\}$
  - $E(G1 \times G2) = \{(e1, e2): e1, e2: \text{ same labels, } p(e1) \text{ and } p(e2) \text{ same labels, } n(e1) \text{ and } n(e2) \text{ same labels}\}$





# Strategies of Design

---

- Kernel as a way to encode prior information
  - Invariance: synonymy, document length, ...
  - Linguistic processing: word normalisation, semantics, stopwords, weighting scheme, ...
- Convolution Kernels: text is a recursively-defined data structure. How to build “global” kernels from local (atomic level) kernels?
- **Generative model-based kernels: the “topology” of the problem will be translated into a kernel function**

# Plan Global

---

- Introduction :
  - Fouille de textes
  - Spécificité des données textuelles
- Approche numéro 1 : méthodes à noyaux
  - Philosophie des méthodes à noyaux
  - Noyaux pour les données textuelles
- Approche numéro 2 : modèles génératifs
  - Génératif versus discriminatif – semi-supervisé
  - Modèles graphiques à variables latentes
  - Exemples : NB, PLSA, LDA, HPLSA
- Perspectives « récentes »

# Generative vs Discriminative

---

## ■ Generative approach:

- Model  $P(x,y)$  ( $= P(y|x).P(x) = P(x|y).P(y)$ )
- Then, for a new  $x$ , choose  $y = \operatorname{argmax} P(x,y)$

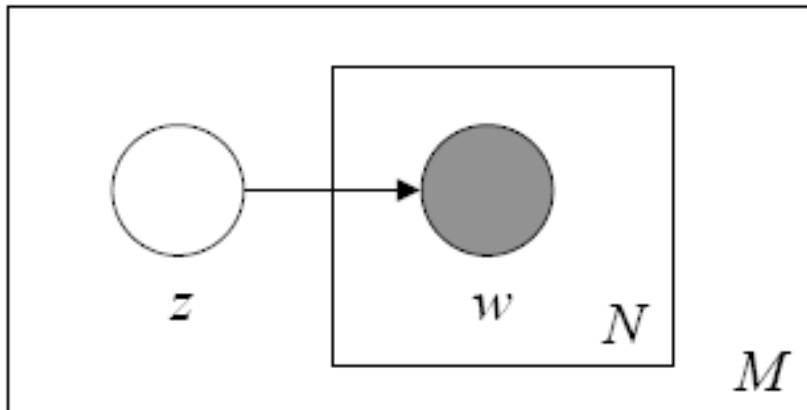
## ■ Discriminative approach:

- Model  $P(y|x)$
- Then, for a new  $x$ , choose  $y = \operatorname{argmax} P(y|x)$

## ■ Most advantages for discriminative approach but:

- Semi-supervised learning – continuum between clustering and categorization
- Novelty Detection
- NB. Most generative approaches use latent variables (hidden classes or components) – strong link between component and categories – Then use probabilistic values of these latent variables as new features in a discriminative setting (cfr. Dimension reduction – generative model-based kernels)

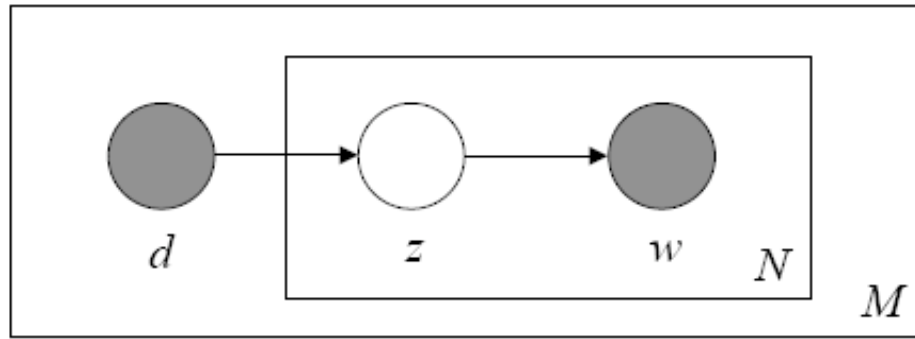
# Graphical models : NB



$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z).$$

- M documents
- N words
- 1! Topic per document
- Supervised case ( $z$  observed):
  - Training : Parameters (class priors and class profiles) by max likelihood
  - Classify :  $\max p(\mathbf{w}, z)$
- Unsupervised:
  - Use EM

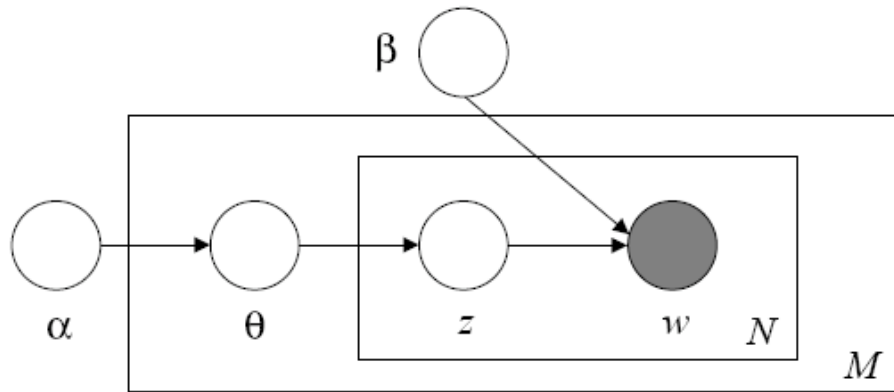
# PLSA



$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d).$$

- M documents
- N words
- Multiple Topics per document
- Supervised case
  - Parameters ( $p(z,d)$  and class profiles) by max likelihood
  - Inference : by EM to identify  $p(z|d)$
- Unsupervised:
  - Use tempered-EM

# LDA



$$p(\mathbf{w}|\alpha, \beta) = \int p(\theta|\alpha) \left( \prod_{n=1}^N \sum_{z_n} p(z_n|\theta) p(w_n|z_n, \beta) \right) d\theta.$$

- M documents
- N words
- Multiple Topics per document
- Dirichlet prior on the topic mixing proportion
- Supervised case
  - Parameters  $(\alpha, \beta)$  (class priors and class profiles) by max likelihood, given  $w, \theta, z$
  - Variational Inference : to identify  $p(\theta, z | \alpha, \beta, \mathbf{w})$
- Unsupervised:
  - Use variational-EM to identify  $(\alpha, \beta)$  , given observed  $\mathbf{w}$

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

# Strategies of Design

---

- Kernel as a way to encode prior information
  - Invariance: synonymy, document length, ...
  - Linguistic processing: word normalisation, semantics, stopwords, weighting scheme, ...
- Convolution Kernels: text is a recursively-defined data structure. How to build “global” kernels from local (atomic level) kernels?
- **Generative model-based kernels: the “topology” of the problem will be translated into a kernel function**



# Remind

---

- This family of strategies brings you the additional advantage of using all your unlabeled training data to design more problem-adapted kernels
- They constitute a natural and elegant way of solving semi-supervised problems (mix of labelled and unlabelled data)

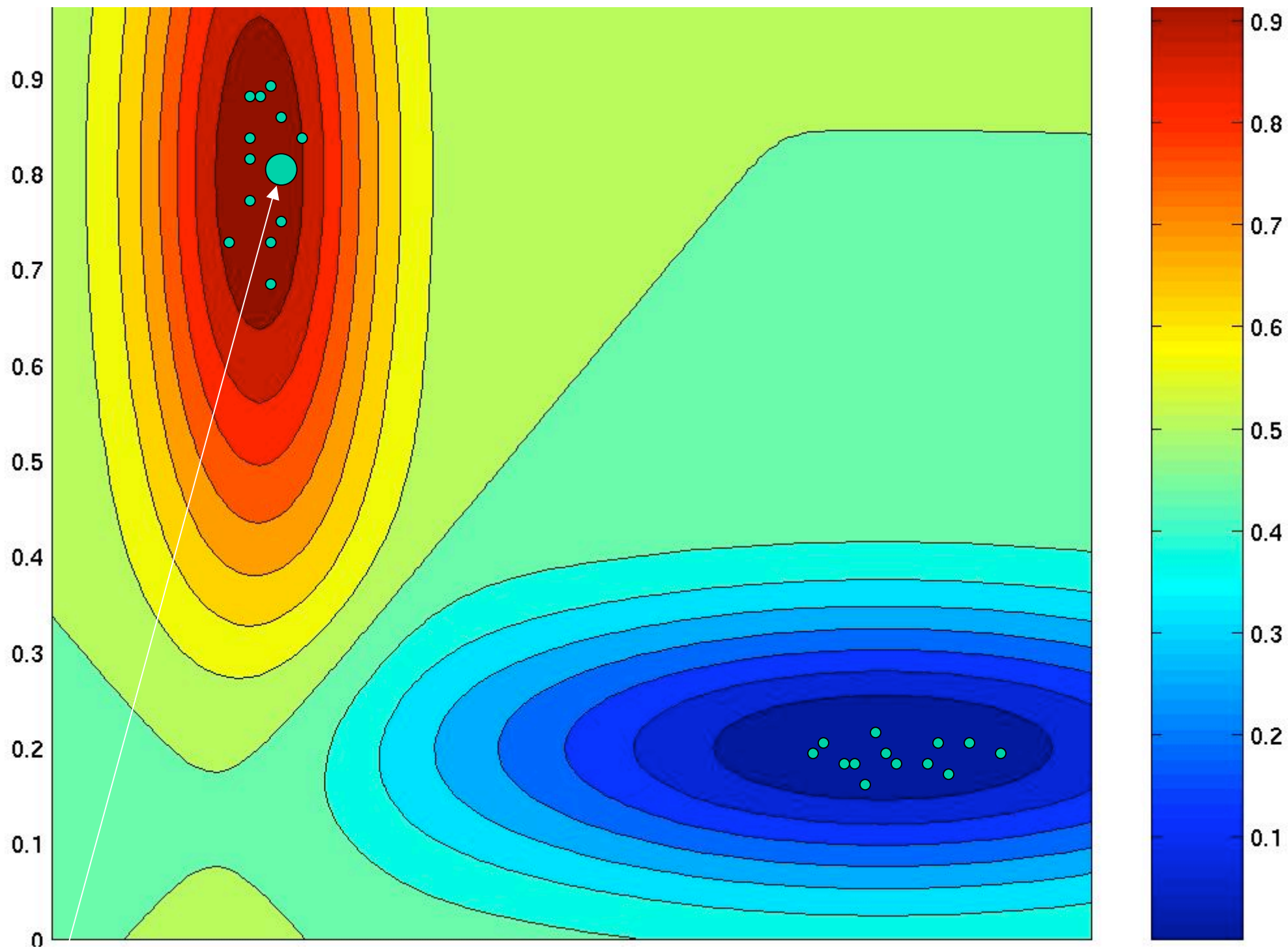
# Marginalised – Conditional Independence Kernels

---

- Assume a family of models  $M$  (with prior  $p_0(m)$  on each model) [finite or countably infinite]
- each model  $m$  gives  $P(x|m)$
- Feature space indexed by models:  $x \rightarrow P(x|m)$
- Then, assuming conditional independence, the joint probability is given by

$$P_M(x, z) = \sum_{m \in M} P(x, z | m) P_0(m) = \sum_{m \in M} P(x | m) P(z | m) P_0(m)$$

- This defines a valid probability-kernel (CI implies PD kernel), by marginalising over  $m$ . Indeed, the gram matrix is  $K = P \cdot \text{diag}(P_0) \cdot P'$  (some reminiscence of latent concept kernels)



# Fisher Kernels

---

- Assume you have only 1 model
  - Marginalised kernel give you little information: only one feature:  $P(x|m)$
  - To exploit much, the model must be “flexible”, so that we can measure how it adapts to individual items → we require a “smoothly” parametrised model
  - Link with previous approach: locally perturbed models constitute our family of models, but  $\dim F = \text{number of parameters}$
- More formally, let  $P(x|\theta_0)$  be the generative model ( $\theta_0$  is typically found by max likelihood); the gradient reflects how the model will be changed to accommodate the new point  $x$  (NB. In practice the loglikelihood is used)  $\nabla_{\dot{\theta}} \log P(x | \dot{\theta}) \Big|_{\dot{\theta}=\dot{\theta}_0}$

# Fisher Kernel : formally

---

- Two objects are similar if they require similar adaptation of the parameters or, in other words, if they stretch the model in the same direction:

$K(x,y)=$

$$(\nabla_{\mathbf{\theta}} \log P(x | \mathbf{\theta})|_{\mathbf{\theta}=\mathbf{\theta}_0})' I_M^{-1} (\nabla_{\mathbf{\theta}} \log P(y | \mathbf{\theta})|_{\mathbf{\theta}=\mathbf{\theta}_0})$$

Where  $I_M$  is the Fisher Information Matrix

$$I_M = E[(\nabla_{\mathbf{\theta}} \log P(x | \mathbf{\theta})|_{\mathbf{\theta}=\mathbf{\theta}_0})(\nabla_{\mathbf{\theta}} \log P(x | \mathbf{\theta})|_{\mathbf{\theta}=\mathbf{\theta}_0})']$$

## Example 2 : PLSA-Fisher Kernels

---

- An example : Fisher kernel for PLSA improves the standard BOW kernel

$$K(d_1, d_2) = \sum_c \frac{P(c|d_1).P(c|d_2)}{P(c)} + \sum_w \tilde{t}f(w, d_1)\tilde{t}f(w, d_2) \sum_c \frac{P(c|d_1, w).P(c|d_2, w)}{P(w|c)}$$

- where  $k_1(d_1, d_2)$  is a measure of how much  $d_1$  and  $d_2$  share the same latent concepts (synonymy is taken into account)
- where  $k_2(d_1, d_2)$  is the traditional inner product of common term frequencies, but weighted by the degree to which these terms belong to the same latent concept (polysemy is taken into account)

# “New” perspectives

---

- Multi-lingual
- Multi-media
- Emotion mining
- Structured documents
- Help to labelling – Active learning