

# Apprentissage Statistique

## Partie 2

Sorbonne Université- Master DAC et Master Math Spécialité Big Data - P. Gallinari,  
[patrick.gallinari@lip6.fr](mailto:patrick.gallinari@lip6.fr), <http://www-connex.lip6.fr/~gallinar/>

nom.prenom@lip6.fr

Année 2018-2019

# Apprentissage supervisé

Machines à noyaux

## Introduction

- ▶ Familles de machines d'apprentissage générales qui exploitent l'idée suivante :
  - ▶ Projeter les données dans un espace de grande dimension - éventuellement infini - où le problème sera facile à traiter
  - ▶ Utiliser des "projections" non linéaires permettant des calculs "efficaces"
- ▶ Exemples :
  - ▶ Machines à Vecteurs Support (généralisent : hyperplan optimal, cadre Vapnik)
  - ▶ Processus Gaussien (généralisent : régression logistique, cadre Bayésien)
- ▶ Cours
  - ▶ Machines à Vecteurs Support

# Perceptron : formulation duale

hyp: 2 classes linéairement séparables, sorties désirées -1, 1

## Perceptron primal

Initialiser  $w(0) = 0$

Répéter (t)

choisir un exemple,  $(x(t), d(t))$

si  $d(t)w(t) \cdot x(t) \leq 0$

alors  $w(t+1) = w(t) + d(t)x(t)$

Jusqu'à convergence

## Fonction de décision

$$F(x) = \operatorname{sgn}\left(\sum_{j=0}^n w_j x_j\right),$$

$$w = \sum_{i=1}^N \alpha_i d^i x^i$$

$\alpha_i$  : nombre de fois où l'algorithme a rencontré une erreur de classification sur  $x^i$

## Perceptron dual

Initialiser  $\alpha = 0, \alpha \in \mathbb{R}^N$

Répéter (t)

choisir un exemple,  $(x(t), d(t))$

soit  $k: x(t) = x^k$

si  $d(t) \sum_{i=1}^N \alpha_i d^i x^i \cdot x(t) \leq 0$

alors  $\alpha_k = \alpha_k + 1$

Jusqu'à convergence

## Fonction de décision

$$F(x) = \operatorname{sgn}\left(\sum_{i=1}^N \alpha_i d^i x^i \cdot x\right)$$

Matrice de Gram  $G$  :

matrice  $N \times N$  de terme  $i, j : x^i \cdot x^j$

matrice de similarité entre données

## Représentation Duale

- ▶ La plupart des machines à apprentissage linéaires ont une représentation duale
  - ▶ Exemples
    - ▶ Adaline, regression, regression ridge, etc
  - ▶ L'information sur les données est entièrement fournie par la matrice de Gram :  $G = (x^i \cdot x^j)_{i,j=1\dots N}$ , qui joue un rôle central
  - ▶ La fonction de décision  $F(x)$  s'exprime comme une combinaison linéaire de produits scalaires entre la donnée d'entrée  $x$  et les exemples d'apprentissage
- ▶ Les machines à noyau généralisent ces idées
  - ▶ Une **fonction noyau**  $K$  est définie sur  $R^n \times R^n$  (on suppose  $x \in R^n$ ) par
$$K(x, z) = \langle \Phi(x), \Phi(z) \rangle$$
où  $\Phi$  est une fonction de  $R^n$  dans un espace muni d'un produit scalaire

# Produit Scalaire et Noyaux

- ▶ **Projection non linéaire dans un espace de grande dimension  $H$** 
  - ▶ (en général  $\dim H \gg n$ , et peut même être infinie)
  - ▶  $\Phi: R^n \rightarrow R^p$  avec  $p \gg n$
- ▶ **Machine linéaire dans  $H$  - Primal :**
  - ▶  $F(x) = \sum_{j=1}^p w_j \Phi_j(x) + b$
- ▶ **Machine linéaire dans  $H$  - Dual :**
  - ▶  $w = \sum_{i=1}^N d^i \alpha_i \Phi(x^i), \quad F(x) = \sum_{i=1}^N d^i \alpha_i \Phi(x^i) \cdot \Phi(x) + b,$
- ▶ **Calculer les produits scalaires dans l'espace initial : choisir  $F /$** 
  - ▶  $\Phi(x) \cdot \Phi(x') = K(x, x')$
  - ▶  $F(x) = \sum_{i=1}^N d^i \alpha_i K(x^i, x) + b$
- ▶ avec  $K$  : fonction noyau (i.e. symétrique)

- ▶ Généralise le produit scalaire dans l'espace initial
- ▶ Le calcul de  $F$  ne dépend pas directement de la taille de  $H$  : les calculs sont faits dans l'espace initial.
- ▶ La machine linéaire dans  $H$  peut être construite à partir d'une fonction  $K$  sans qu'il soit nécessaire de définir explicitement  $\Phi$  : en pratique, on spécifiera directement  $K$ .
- ▶ Cette idée permet d'étendre de nombreuses techniques linéaires au non linéaire: il suffit de trouver des noyaux appropriés
  - ▶ Exemples
    - ▶ ACP, analyse discriminante, regression, etc

## Caractérisation des noyaux

▶ Quand peut on utiliser cette idée ?

▶ Cas d'un espace fini

- ▶ Soit  $X = \{x^1, \dots, x^N\}$ ,  $K(x, x')$  une fonction symétrique sur  $X$ ,  $K$  est une fonction noyau ssi la matrice  $N \times N$  dont l'élément  $(i, j)$  est  $K(x^i, x^j)$  est positive semi-définie (valeurs propres  $\geq 0$  ou  $x^T K x > 0 \forall x$ , avec  $K = \text{Matrice} (K(x^i, x^j)); i, j = 1..N$ )

▶ Cas général : Conditions de Mercer (noyaux de Mercer)

- ▶ Il existe une application  $\Phi$  et un développement

$$K(x, x') = \sum_{i=1}^{\infty} \Phi(x)_i \cdot \Phi(x')_i$$

- ▶ Ssi :

$$\forall g / \int g(x)^2 dx \text{ est fini et } \int K(x, x') g(x) g(x') dx dx' \geq 0$$



# Caractérisation des noyaux

## Espace de Hilbert à noyau autoreproduisant

- ▶ Une fonction  $K: X * X \rightarrow R$  qui est soit continue soit définie sur un domaine fini peut s'écrire sous la forme d'un produit scalaire :

$$K(x, z) = \langle \phi(x), \phi(z) \rangle$$

avec  $\Phi : x \rightarrow \Phi(x) \in \mathbf{F}$  espace de Hilbert

ssi c'est une fonction symétrique et toutes les matrices formées par la restriction de  $K$  à un échantillon fini sur  $X$  sont semi-définies positives).

- ▶ Résultat à la base de la caractérisation effective des fonctions noyaux
- ▶ Il permet de caractériser  $K$  comme un noyau sans passer par  $\Phi$
- ▶ C'est une formulation équivalente aux conditions de Mercer

# Caractérisation des noyaux

## Espace de Hilbert à noyau autoreproduisant

- ▶ L'espace de Hilbert associé au noyau  $K$  :

$$F = \left\{ \sum_{i=1}^l \alpha_i K(x_i, \cdot) \mid l \in \mathbb{N}, x_i \in X, \alpha_i \in \mathbb{R}, i = 1..l \right\}$$

- ▶ Le produit scalaire défini sur cet espace :

$$\begin{aligned} \text{Soient } f(\cdot) &= \sum_{i=1}^l \alpha_i K(x_i, \cdot), & g(\cdot) &= \sum_{j=1}^n \beta_j K(x_j, \cdot) \\ \langle f, g \rangle &= \sum_{i=1}^l \sum_{j=1}^n \alpha_i \beta_j K(x_i, z_j) = \sum_{i=1}^l \alpha_i g(x_i) = \sum_{j=1}^n \beta_j f(z_j) \end{aligned}$$

▶ Noyau auto-reproduisant

▶ Si on prend  $g(\cdot) = K(x, \cdot)$  alors  $\langle f, K(x, \cdot) \rangle = \sum_{i=1}^l \alpha_i K(x_i, x) = f(x)$

## Exemples de noyaux

$$K(x, z) = \langle x \cdot z \rangle^2$$

$$K(x, z) = \left( \sum_{i=1}^n x_i \cdot z_i \right)^2 = \sum_{i,j=1}^n (x_i \cdot x_j)(z_i \cdot z_j)$$

$$K(x, z) = \langle \phi(x) \cdot \phi(z) \rangle \text{ avec } \phi(x) / \phi(x)_{i,j} = (x_i \cdot x_j)_{i,j=1,n}$$

i.e. tous les monomes de d° 2:  $\binom{n+1}{2}$

$$K(x, z) = (\langle x \cdot z \rangle + c)^2$$

$$K(x, z) = \langle \phi(x) \cdot \phi(z) \rangle \text{ avec } \phi(x) / \phi(x)_{i,j} = ((x_i \cdot x_j)_{i,j=1,n}, (\sqrt{2c}x_i)_{i=1,n}, c)$$

i.e. ss ensemble des polynomes de d° 2

## Exemples de noyaux (suite)

$$K(x, x^i) = \begin{cases} (x \cdot x^i + 1)^d & \Phi \text{ polynome d'ordre } d \\ \exp - \gamma \|x - x^i\|^2 & \Phi \text{ noyau gaussien} \\ \textit{Sigmoide}(vx \cdot x^i + c) & \end{cases}$$

- ▶ En pratique, on utilise souvent des noyaux
  - ▶ Linéaires
  - ▶ Gaussiens
  - ▶ Polynomiaux

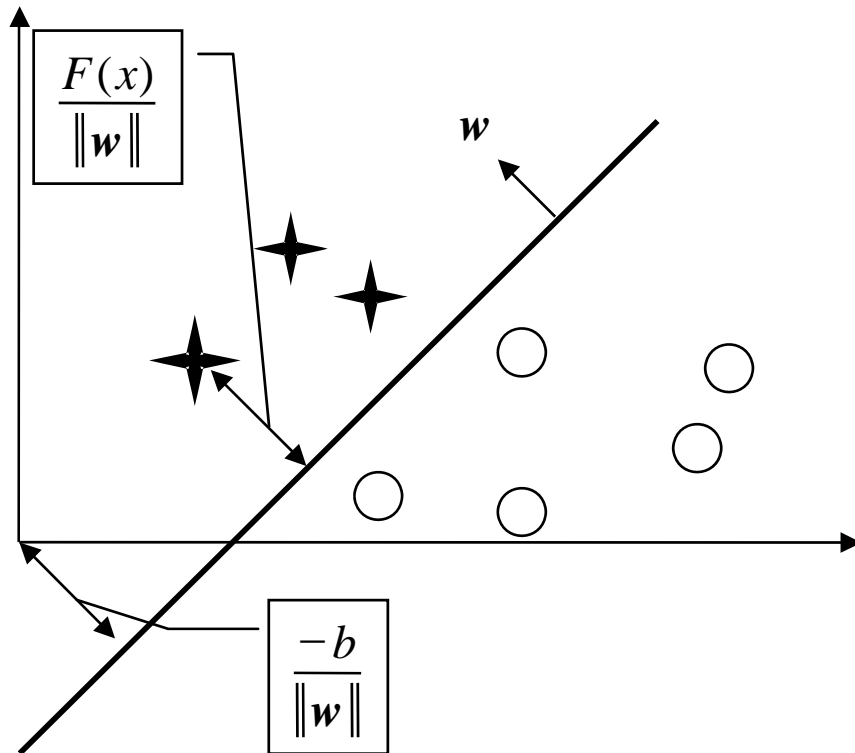
## Construction des noyaux en pratique

- ▶ Les résultats de Mercer servent à prouver les propriétés des fonctions noyaux. En pratique, elles sont peu utiles
- ▶ Pour construire des noyaux, on procède par combinaison à partir de noyaux connus
- ▶ Si  $K_1$  et  $K_2$  sont des noyaux sur  $X^2$ ,  $K_3$  défini sur  $F$ , les fonctions suivantes sont des noyaux :
  - ▶  $K(x, z) = K_1(x, z) + K_2(x, z)$
  - ▶  $K(x, z) = K_1(x, z) \cdot K_2(x, z)$
  - ▶  $K(x, z) = aK_1(x, z)$
  - ▶  $K(x, z) = K_3(\phi(x), \phi(z))$
  - ▶ .....

## Machines à vecteurs support

- ▶ Exposé du cours : discrimination 2 classes
- ▶ Cas général : discrimination multi-classes, régression, densité
- ▶ Idées
  - ▶ Projeter -non linéairement- les données dans un espace de "très" grande taille  $H$
  - ▶ Faire une séparation linéaire de bonne qualité dans cet espace
  - ▶ Raisonner dans  $H$ , mais résoudre le problème d'optimisation dans l'espace de départ (noyaux)

► Notion de marge



► Hyperplan H

►  $F(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = 0$

► Marge géométrique pour  $x^i$

►  $M(x^i) = d^i \left( \frac{\mathbf{w} \cdot \mathbf{x}^i}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|} \right) = d^i \frac{F(x^i)}{\|\mathbf{w}\|}$

► Marge de  $\mathbf{w}$  par rapport à un ensemble de données  $D$

►  $\text{Min}_i M(x^i)$

► Hyperplan de marge maximale

►  $\text{Max}_{\mathbf{w}} (\text{Min}_i M(x^i))$



## Marge géométrique vs marge fonctionnelle

- ▶ Marge géométrique

- ▶  $d^i \frac{F(x^i)}{\|w\|}$

- ▶ Marge fonctionnelle

- ▶  $d^i F(x^i)$

- ▶ Remplacer  $w$  par  $kw$  ne change pas la fonction de décision ou la marge géométrique, mais change la marge fonctionnelle.

- ▶ Pour les SVM, on fixera la marge fonctionnelle à 1 et on optimisera la marge géométrique.

## Prémises : Séparation linéaire à hyperplan optimal (1974)

▶ Hyp :  $D$  linéairement séparable

▶ Fonction de décision :  $F(x) = w \cdot x + b$   $D = \{x^i, d^i\}_{i=1..N}$  avec  $d^i = \pm 1$

▶ Pb apprentissage :

▶ trouver l'hyperplan optimal  $H^*$  qui sépare  $D$  i.e.

▶  $d^i \cdot F(x^i) \geq 1, \forall i$

▶ avec une marge géométrique maximale  $M = \min_i \frac{d^i \cdot F(x^i)}{\|w\|} = \frac{1}{\|w\|}$

i.e. : Problème Primal :

$$\begin{cases} \text{Minimiser } \|w\|^2 \\ \text{S.C. } d^i \cdot F(x^i) \geq 1 \end{cases}$$

- 
- ▶ Dans le cas de noyaux linéaires, on résout en général le problème primal
    - ▶ Il existe de nombreux algorithmes rapides
  - ▶ Dans le cas de noyaux non linéaires, on va en général résoudre une version duale du problème
    - ▶ On introduit ensuite quelques notions d'optimisation sous contrainte pour décrire la formulation duale du problème

# Intermède

Optimisation  
Problèmes sous contraintes égalités, inégalités

# Optimisation

## Problèmes sous contraintes égalités, inégalités

### ▶ Soient

- ▶  $f, g_{i,i=1,\dots,k}, h_{j,j=1,\dots,m}$  des fonctions définies sur  $\mathbb{R}^n$  à valeur dans  $\mathbb{R}$

### ▶ On considère le problème primal suivant (pb. (0)) :

$$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$$

Sous contraintes

$$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k$$

$$\text{noté } \mathbf{g}(\mathbf{w}) \leq 0$$

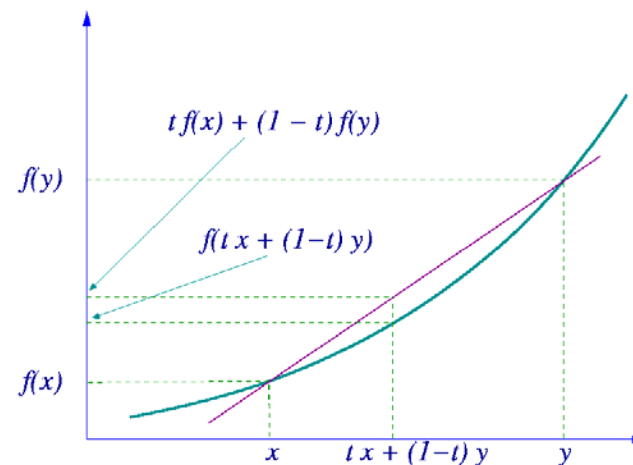
$$h_j(\mathbf{w}) = 0, j = 1, \dots, m$$

$$\text{noté } \mathbf{h}(\mathbf{w}) = 0$$

- ▶ **Fonction objectif**  $f(\mathbf{w})$
- ▶ **Région admissible**  $R = \{\mathbf{w} \in \Omega: \mathbf{g}(\mathbf{w}) \leq 0, \mathbf{h}(\mathbf{w}) = 0\}$ , région de  $\Omega$  où  $f$  est définie et les contraintes vérifiées
- ▶  $\mathbf{w}^*$  est un **minimum global** si il n'existe pas d'autre point  $\mathbf{w}$  tel que  $f(\mathbf{w}) < f(\mathbf{w}^*)$ , c'est un optimum local si  $\exists \epsilon > 0: f(\mathbf{w}) \geq f(\mathbf{w}^*)$ , sur la boule  $\|\mathbf{w} - \mathbf{w}^*\| < \epsilon$
- ▶ Une contrainte  $g_i(\mathbf{w}) \leq 0$  est dite **active** si la solution  $\mathbf{w}^*$  vérifie  $g_i(\mathbf{w}^*) = 0$  et inactive sinon
- ▶ La valeur optimale de la fonction objectif (solution du pb. (0)) est appelée la **valeur du problème d'optimisation** primal)

## Optimisation - Fonctions convexes

- ▶  $f(w)$  est convexe pour  $w \in \mathbb{R}^n$  si
- ▶  $\forall t \in [0,1], \forall w, v \in \mathbb{R}^n, \forall t \in [0,1], f(tw + (1-t)v) \leq tf(w) + (1-t)f(v)$



- ▶ Un ensemble  $\Omega \subset \mathbb{R}^n$  est convexe si  $\forall w, v \in \mathbb{R}^n, \forall t \in [0,1], tw + (1-t)v \in \Omega$
- ▶ Si une fonction est convexe, tout minimum local est un minimum global
- ▶ Un problème d'optimisation pour lequel  $\Omega$  est convexe, la fonction objectif et les contraintes sont convexes est dit convexe

## Optimisation non contrainte

### ▶ Th. Fermat

- ▶ Une Condition Nécessaire pour que  $w^*$  soit un min. de  $f(w)$ ,  $f \in C^1$  est  $\frac{\partial f(w^*)}{\partial w} = 0$
- ▶ Si  $f$  est convexe c'est une Condition Suffisante

## Optimisation avec contraintes égalités Lagrangien

- ▶ Optimisation avec contraintes égalité (pb (I)):

$$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$$

Sous les contraintes

$$h_j(\mathbf{w}) = 0, \quad j = 1, \dots, m \quad \text{noté } \mathbf{h}(\mathbf{w}) = 0$$

- ▶ On définit le Lagrangien  $L(\mathbf{w}, \boldsymbol{\beta})$  associé à ce problème par

$$L(\mathbf{w}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{j=1}^m \beta_j h_j(\mathbf{w})$$

- ▶ les  $\beta_j$  sont les coefficients de Lagrange
- ▶ Rq.
  - ▶ Si  $\mathbf{w}^*$  est une solution du problème d'optimisation sous contrainte, Il est possible que  $\frac{\partial f(\mathbf{w}^*)}{\partial \mathbf{w}} \neq 0$



# Optimisation avec contraintes égalités

## Th. Lagrange

### ▶ Th. Lagrange

- ▶ Une condition nécessaire pour que  $\mathbf{w}^*$ , soit solution de (pb. (I)), avec  $f, h_i \in C^1$  est

- $\frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \mathbf{w}} = 0$

- $\frac{\partial L(\mathbf{w}^*, \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}} = 0$

- ▶ Si  $L(\mathbf{w}, \boldsymbol{\beta}^*)$  est une fonction convexe de  $\mathbf{w}$ , c'est une condition suffisante

### ▶ Rq

- ▶ La première condition donne un nouveau système d'équations
  - ▶ La seconde donne les contraintes

## Optimisation sous contraintes égalité + inégalités - Lagrangien augmenté

- ▶ De même, on définit le Lagrangien augmenté pour le pb. (0) :

$$L(\mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = f(\mathbf{w}) + \sum_{i=1}^k \alpha_i g_i(\mathbf{w}) + \sum_{j=1}^m \beta_j h_j(\mathbf{w})$$

- ▶ Rappel : pb. (0)

*Min*  $f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$

Sous contraintes

$$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k$$

$$h_j(\mathbf{w}) = 0, j = 1, \dots, m$$

$$\text{noté } \mathbf{g}(\mathbf{w}) \leq 0$$

$$\text{noté } \mathbf{h}(\mathbf{w}) = 0$$

## Optimisation – formulation duale

### ▶ Formulation duale du problème d'optimisation

- ▶ Le problème d'optimisation dual correspondant au problème primal pb (0) est :

$$\text{Maximiser}_{\alpha, \beta} \theta(\alpha, \beta) = \min_{w \in \Omega} L(w, \alpha, \beta)$$

sous contrainte  $\alpha \geq 0$

- ▶ Max  $\theta(\alpha, \beta)$  est appelé la **valeur du dual**
- ▶ Rq :  $\min_{w \in \Omega} L(w, \alpha, \beta)$  est une fonction de  $\alpha, \beta$  uniquement

- ▶ Propriété (dualité faible) : la valeur du dual est bornée supérieurement par la valeur du primal

$$\max_{\alpha, \beta; \alpha \geq 0} \min_{w \in \Omega} L(w, \alpha, \beta) \leq \min_{w \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta)$$

Dans certains cas, on a égalité, cf. dualité forte

## Optimisation : dualité faible

- ▶  $\max_{\alpha, \beta; \alpha \geq 0} \min_{v \in \Omega} L(v, \alpha, \beta) \leq \min_{w \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta)$ :
  - ▶  $\theta_D(\alpha, \beta) = \min_{v \in \Omega} L(v, \alpha, \beta)$ :
    - ▶  $\theta_D(\alpha, \beta) \leq L(w, \alpha, \beta)$  pour un  $w$  qqe admissible
    - ▶  $L(w, \alpha, \beta) = f(w) + \alpha \cdot g(w) + \beta \cdot h(w)$  avec  $\alpha \geq 0, g(w) \leq 0, h(w) = 0$
    - ▶  $L(w, \alpha, \beta) \leq f(w)$
    - ▶  $\theta_D(\alpha, \beta) \leq f(w)$
    - ▶  $\max_{\alpha, \beta; \alpha \geq 0} \min_{v \in \Omega} L(v, \alpha, \beta) \leq f(w)$  pour un  $w$  qqe admissible
  - ▶  $\theta_P(w) = \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta)$ 
    - ▶  $\theta_P(w) = f(w) + \max_{\alpha, \beta; \alpha \geq 0} (\alpha \cdot g(w) + \beta \cdot h(w))$
    - ▶  $\theta_P(w) = \begin{cases} f(w) & \text{si } w \text{ est admissible} \\ +\infty & \text{sinon} \end{cases}$  car  $\max_{\alpha, \beta; \alpha \geq 0} (\alpha g(w) + \beta h(w)) = 0$  pour un point admissible
    - ▶  $\theta_P(w) = f(w)$
    - ▶ Rq: on retrouve le problème primal original :  $\min_{w \in \Omega} \max_{\alpha, \beta; \alpha \geq 0} L(w, \alpha, \beta) = \min_{w \in \Omega} f(w)$
- ▶ D'où l'inégalité de dualité faible

## ▶ Théorème de dualité forte

### ▶ Etant donné un problème d'optimisation

$\text{Min } f(\mathbf{w}), \mathbf{w} \in \Omega \subset \mathbb{R}^n$  convexe et  $f \in \mathcal{C}^1$  convexe

Sous contraintes

$$g_i(\mathbf{w}) \leq 0, i = 1, \dots, k \quad \text{noté } \mathbf{g}(\mathbf{w}) \leq 0$$

$$h_j(\mathbf{w}) = 0, j = 1, \dots, m \quad \text{noté } \mathbf{h}(\mathbf{w}) = 0$$

où les  $g_i$  et les  $h_j$  sont affines ( $h_j(\mathbf{w}) = A_j\mathbf{w} + b_j$ )

- ▶ alors les valeurs du primal et du dual sont égales
- ▶ Les conditions d'existence d'un optimum sont données par le théorème de Kuhn et Tucker

# Optimisation

## Th. Kuhn et Tucker

- ▶ On considère (pb. (0)) avec  $\Omega$  convexe et  $f \in C_1$  convexe,  $g_i, h_j$  affines ( $h = A.w + b$ )

- ▶ Une CNS pour que  $w^*$  soit un optimum est qu'il existe  $\alpha^*$  et  $\beta^*$  :

$$\left\{ \begin{array}{l} \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial w} = 0 \\ \frac{\partial L(w^*, \alpha^*, \beta^*)}{\partial \beta} = 0 \\ \alpha_i^* g_i(w^*) = 0, i = 1..k \\ g_i(w^*) \leq 0, i = 1..k \\ \alpha_i^* \geq 0, i = 1..k \end{array} \right.$$

- ▶ Sous les hypothèses de convexité, la formulation duale est une alternative à la formulation primale qui peut se révéler plus simple à traiter (e.g. SVM non linéaires)

## Optimisation

### ► Rq

- La 3<sup>e</sup> condition dite condition complémentaire de Karush-Kuhn-Tucker implique que pour une contrainte active  $\alpha_i^* \geq 0$  alors que pour une contrainte inactive  $\alpha_i^* = 0$ 
  - Soit une contrainte est **active** ( $\alpha_i^* \geq 0$  et  $g_i(w^*) = 0$ ),  $w^*$  est un point frontière de la région admissible
  - Soit elle est **inactive** ( $\alpha_i^* = 0$ ) et  $w^*$  est dans la région admissible
- Si le point solution  $w^*$  est dans la région admissible (contrainte inactive) alors les conditions d'optimalité sont données par le th. de Fermat et  $\alpha_i^* = 0$ . Si il est sur la frontière (contrainte active), les conditions d'optimalité sont données par le th. de Lagrange avec  $\alpha_i^* > 0$ .

### ► Fin de l'intermède

# SVM – formulations primale et duale

Cas d'un noyau linéaire

- ▶ SVM

- ▶  $\Omega, f$  et les contraintes sont convexes,  $L$  est quadratique
- ▶ On étudie le cas,  $D = \{(x^i, d^i)\}_{i=1\dots N}$  linéairement séparables,  $d^i \in \{-1, 1\}$

- ▶ Pb. Primal  $Min(w.w)$  (i.e. max la marge)

sous les contraintes

$$d^i (w.x^i + b) \geq 1, i = 1..N$$

- ▶ Lagrangien primal

$$L(w, b, \alpha) = \frac{1}{2} w.w - \sum_{i=1}^N \alpha_i (d^i (w.x^i + b) - 1)$$

- ▶ Lagrangien dual

$$L(w, b, \alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N d^i d^j \alpha_i \alpha_j (x^i . x^j)$$

- ▶ Avec  $\alpha_i \geq 0$  dans les 2 cas



# SVM – formulations primale et duale

## ▶ Pb. Dual

$$\text{Max } L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N d^i d^j \alpha_i \alpha_j (x^i \cdot x^j)$$

sous les contraintes

$$\begin{cases} \sum_{i=1}^N d^i \alpha_i = 0 \\ \alpha_i \geq 0, i = 1..N \end{cases}$$

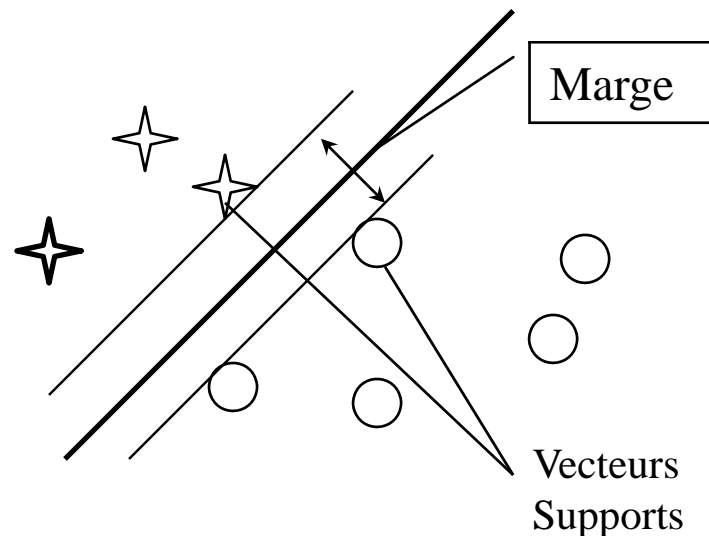
## ▶ C'est un problème d'optimisation quadratique sous contrainte

- ▶ Solution :  $w^*$  dépend uniquement des points supports, i.e. des points sur la marge qui vérifient :  $d^i F^*(x^i) = 1$

- ▶ la fonction de décision prend la forme

$$F(x, \alpha^*, \beta^*) = \sum_{\text{vecteurs support}} d^i \alpha^*_i (x^i \cdot x) + b^*$$

- ▶ Rq: Quelque soit la dimension de l'espace, le nombre de degrés de liberté est "égal" au nombre de points de support
- ▶  $F^*$  dépend uniquement du produit scalaire  $x^i \cdot x$



## Machines à vecteurs supports cas de noyaux non linéaires

- ▶ Faire une séparation à marge max. dans un espace défini par une fonction noyau.
- ▶ Tous les résultats sur le classifieur linéaire à marge max. se transposent en remplaçant  $x^i \cdot x$  par  $K(x^i, x)$

$$\Phi : R^n \rightarrow R^p$$
$$W = \sum_{x^i \in V.S.} d^i \alpha_i \Phi(x^i) \qquad F(x) = \sum_{x^i \in V.S.} d^i \alpha_i \Phi(x^i) \Phi(x) + b$$

$$\Phi(x) \cdot \Phi(x') = K(x, x')$$
$$F(x) = \sum_{x^i \in V.S.} d^i \alpha_i K(x, x^i) + b$$

## Machines à vecteurs supports

- ▶ Apprentissage :

- ▶ On résout le problème d'optimisation dual :

$$\text{Maximiser } L(\alpha) = \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j d^i d^j K(x^i, x^j)$$

$$\text{S.C } \alpha_i \geq 0 \text{ et } \sum_i \alpha_i d^i = 0$$

- ▶ Problème minimisation quadratique sous contraintes dans l'espace de départ
- ▶ Difficile en pratique : différents algorithmes.
- ▶ Dans la solution optimale  $\alpha_i > 0$  uniquement pour les points support.
  - ▶ Seuls les produits scalaires  $K$  apparaissent, et pas les  $\Phi$ .

# Propriétés de généralisation -exemples

- ▶ Th 1  $E [P (erreur(x))] \leq \frac{E [\# \text{ vecteurs supports}]}{\# \text{ exemples apprentissage} - 1}$
- ▶ peu de points support  $\rightarrow$  meilleure généralisation
  - ▶ indépendant de la taille de l'espace de départ

▶ Th 2

- ▶ Si  $\exists q / \forall i = 1..N, \|x^i\| \leq q$
- ▶ l'hyperplan optimal passe par l'origine et a pour marge  $\rho$

- ▶ Alors  $E [P (erreur(x))] \leq \frac{E [\frac{q}{2}]}{N}$

- Dans les 2 cas,  $E[P()]$  est l'espérance sur tous les ensembles de taille  $l-1$ , et  $E[\text{membre droit}]$  est l'espérance sur tous les ensembles d'apprentissage de taille  $l$  (leave one out).

# Machines à vecteurs supports

## Cas non linéairement séparable

- ▶ Marges molles
  - ▶ L'algorithme est instable
    - ▶ Dans les cas non linéairement séparables
    - ▶ Dans le cas de données réelles même linéairement séparables
    - ▶ **Solution adoptée en pratique**
      - autoriser des erreurs, i.e. prendre pour contraintes :

$$d^i (W \cdot \Phi(x^i) + b) \geq 1 - \eta^i$$

$$\eta^i \geq 0$$

- ▶  $\eta^i = 0$ ,  $x^i$  est correctement classifié et est du bon coté de la marge
- ▶  $0 < \eta^i \leq 1$ ,  $x^i$  est correctement classifié, est à l'intérieur de la marge
- ▶  $\eta^i > 1$ ,  $x^i$  est mal classé
- ▶  $\eta^i$  : *slack variable*

# Machines à vecteurs supports

## Cas non linéairement séparable

- ▶ But
  - ▶ Maximiser la marge tout en pénalisant les points qui sont mal classés
- ▶ Formalisation
  - ▶ Plusieurs expressions possibles du problème
  - ▶ L'une des plus courantes :

$$\text{Min}(w.w) + C \sum_{i=1}^N \eta^i \quad (\text{i.e. max la marge})$$

S.C.

$$d^i (w.x^i + b) \geq 1 - \eta^i, i = 1..N$$

$$\eta^i \geq 0, i = 1..N$$

- ▶ C fixé par validation croisée joue le rôle de paramètre de régularisation

# Machines à vecteurs supports

## Cas non linéairement séparable

- ▶ Marges molles – formulation duale

$$\text{Maximiser} \quad L(\alpha) = \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j d^i d^j K(x^i, x^j)$$

$$\text{S.C} \quad 0 \leq \alpha_i \leq C \quad \text{et} \quad \sum_i \alpha_i d^i = 0$$



# Algorithmes d'optimisation

- ▶ **Algorithmes d'optimisation standard pour la programmation quadratique sous contrainte**
  - ▶ e.g. **Sequential Minimal Optimization (SMO)**
- ▶ **Algorithmes stochastiques - SVM Results –(Bottou 2007)**
  - ▶ Task : Document classification - RCV1 documents belonging to the class CCAT (2 classes classification task)
    - ▶ Programs [SVMLight](#) and [SVMPerf](#) are well known SVM solvers written by [Thorsten Joachims](#). SVMLight is suitable for SVMs with arbitrary kernels. Similar results could be achieved using [Chih-Jen Lin's LibSVM](#) software. SVMPerf is a specialized solver for linear SVMs. It is considered to be one of the most efficient optimizer for this particular problem.

Algorithm (hinge loss)	Training Time	Primal cost	Test Error
<a href="#">SVMLight</a>	23642 secs	0.2275	6.02%
<a href="#">SVMPerf</a>	66 secs	0.2278	6.03%
Stochastic Gradient (svmsgd)	<b>1.4 secs</b>	0.2275	6.02%
Stochastic Gradient (svmsgd2)	<b>1.4 secs</b>	0.2275	6.01%

# Apprentissage non supervisé

Algorithme EM et mélange de densités

Spectral clustering

Non Negative Matrix Factorization

# Applications

- ▶ analyse des données quand il n'y a pas de connaissance sur la classe.
  - ▶ e.g. pas d'étiquetage des données (problème nouveau)
- ▶ trop de données ou étiquetage trop compliqué
  - ▶ e.g. traces utilisateur (web), documents web, parole, etc
- ▶ réduction de la quantité d'information
  - ▶ e.g. quantification
- ▶ découverte de régularités sur les données ou de similarités.

# Apprentissage non supervisé

Algorithme Espérance Maximisation (EM)  
Application aux mélanges de densités

## Algorithme E. M. (Espérance Maximisation) - Introduction

- ▶ On dispose
  - ▶ de données  $D = \{x^i; i = 1 \dots N\}$ 
    - ▶ On n'a pas d'étiquette  $d^i$  associée à  $x^i$
  - ▶ d'un modèle génératif, de paramètres  $\theta$
- ▶ On veut trouver les paramètres du modèle qui expliquent au mieux la génération des données
- ▶ On se donne un critère
  - ▶ Ici on considère la vraisemblance des données qui est le critère le plus fréquent
    - ▶  $P(D; \theta) = P(x^1, \dots, x^N; \theta)$
  - ▶ D'autres critères sont également utilisés
- ▶ On va essayer de déterminer les paramètres  $\theta$  de façon à maximiser la vraisemblance

## EM Introduction

### Exemple : mélange de deux populations

- ▶ On recueille des données sur deux populations
  - ▶ e.g. taille d'individus  $D = \{x^i; i = 1 \dots N\}$
  - ▶ Hypothèse : les données de chaque population sont gaussiennes à 1 dimension
    - ▶  $N(\mu_1, \sigma_1), N(\mu_2, \sigma_2)$
- ▶ Problème
  - ▶ Estimer les  $\mu_i$  et les  $\sigma_i$  à partir des données
  - ▶ Si les  $d^i$  sont connus, i.e.  $D = \{(x^i, d^i); i = 1 \dots N\}$  la solution est simple
    - ▶ On a deux population séparées (2 classes)  $C_1, C_2$
    - ▶ On utilise les estimateurs classiques de la moyenne et de la variance
      - Pour la moyenne  $\mu_j = \frac{1}{|C_j|} \sum_{x^i \in C_j} x^i$ , idem pour la variance
    - ▶ Ces estimateurs usuels sont les estimateurs du maximum de vraisemblance
      - cf. slide suivant

## EM Introduction - Mélange de deux populations

### Cas où l'appartenance est connue

#### ▶ Vraisemblance

$$▶ P(D|\theta) = \prod_{x^i \in C_1} p(x^i | \theta_1) \prod_{x^j \in C_2} p(x^j | \theta_2)$$

#### ▶ En pratique on maximise la log-vraisemblance

$$▶ L(\theta) = \log(P(D|\theta)) = \sum_{x^i \in C_1} \log p(x^i | \theta_1) + \sum_{x^j \in C_2} \log p(x^j | \theta_2)$$

#### ▶ Cas des gaussiennes

$$▶ p(x|C_k) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma_k^2}\right)$$

$$▶ \frac{\partial L}{\partial \mu_k} = 0 \Leftrightarrow \mu_k = \frac{1}{|C_k|} \sum_{x^i \in C_k} x^i, \text{ idem pour la variance}$$

## EM Introduction - Mélange de deux populations

### Cas où la probabilité d'appartenance est connue

- ▶ On connaît les  $p(C_k|x)$ ,  $k = 1, 2$  pour tous les  $x$  de l'ensemble d'apprentissage
- ▶  $p(x^i|\theta) = p(x^i|C_1)p(C_1) + p(x^i|C_2)p(C_2)$
- ▶ C'est un modèle de mélange de deux densités
- ▶ **Log-vraisemblance**
  - ▶  $L(\theta) = \log(P(D|\theta)) = \sum_{x^i} \log(p(x^i|C_1)p(C_1) + p(x^i|C_2)p(C_2))$
  - ▶ Cas des gaussiennes
    - ▶  $\frac{\partial L}{\partial \mu_k} = 0 \Leftrightarrow \mu_k = \frac{\sum_{x^i} p(C_k|x^i) \cdot x^i}{\sum_{x^i} p(C_k|x^i)}$ , idem pour la variance



## EM Introduction - Mélange de densités gaussiennes

### La probabilité d'appartenance est inconnue

- ▶ On suppose que les données sont générées par le modèle suivant
  - ▶ Tirer  $z^i \sim \text{Multinomiale}(\phi)$   $\phi_i \geq 0, \sum_{j=1}^k \phi_j = 1$
  - ▶ Tirer  $x^i$  /  $x^i | z^i = j \sim N(\mu_j, \sigma_j)$
  - ▶ i.e.  $x^i$  est généré en choisissant d'abord  $z^i \in \{1, \dots, k\}$  et ensuite en effectuant un tirage suivant une gaussienne  $N(\mu_j, \sigma_j)$  si  $z^i = j$
- ▶ C'est ce qu'on appelle un mélange de gaussienne
- ▶ la vraisemblance des données s'écrit
  - ▶  $l(\theta) = l(\mu, \sigma, \phi) = \prod_{i=1}^N p(x^i; \theta) = \prod_{i=1}^N \sum_{j=1}^k p(x^i, z^i = j; \theta)$ 
$$= \prod_{i=1}^N \sum_{j=1}^k p(x^i | z^i = j; \mu, \sigma) p(z^i | \phi)$$
- ▶ La log vraisemblance
  - ▶  $l(\mu, \sigma, \phi) = \sum_{i=1}^N \log(\sum_{j=1}^k p(x^i | z^i = j; \mu, \sigma) p(z^i | \phi))$
- ▶ Les deux exemples précédents (appartenance connue et probabilité d'appartenance connue sont des cas particuliers de ce modèle de mélange)
- ▶ Quand la probabilité d'appartenance est inconnue, on ne peut pas trouver une forme analytique pour les estimateurs du maximum de vraisemblance

## Algorithme E.M. – Cas général

### ▶ Problème

- ▶ On a un problème d'estimation pour lequel on connaît un ensemble d'apprentissage  $D = \{x^1, \dots, x^N\}$
- ▶ On veut trouver les paramètres qui maximisent la (log) vraisemblance des données  $L(\theta) = \log p(D; \theta)$
- ▶ On ne connaît pas de formule analytique permettant d'estimer les paramètres  $\theta$

### ▶ On postule

- ▶ l'existence de variables cachées  $z$  responsables de la génération des données
  - ▶ À chaque  $x^i$ , on associe sa classe cachée  $z^i$   
 $Z = \{z^i; i = 1..N\}$
- ▶ l'existence d'une fonction densité jointe sur les données observées et cachées  $p(x, z)$ 
  - ▶  $p(D, Z|\theta)$  sera appelé **vraisemblance complète** des données pour le modèle  $\theta$ .

### ▶ Remarque

- ▶ Les variables  $z$  sont inconnues et sont considérées comme des variables aléatoires
- ▶  $P(D, Z|\theta)$  sera donc elle-même une variable aléatoire

## Algorithme EM - Cas général

- ▶ L'algorithme E.M. fournit une solution à notre problème
- ▶ C'est un algorithme itératif en 2 étapes
  - ▶ Etape E. (Espérance)
    - ▶ Construire une borne inférieure de  $L(\theta)$
    - ▶ Cette borne inférieure est **l'espérance** de la vraisemblance complète des données
  - ▶ Etape M. (Maximisation)
    - ▶ Déterminer les paramètres, i.e. les  $\phi_j, \mu_j, \sigma_j, j = 1, \dots, k$  qui optimisent cette borne inférieure

## Algorithme EM – Cas general

- ▶ La borne inférieure est définie par une fonction auxiliaire  $Q$ 
  - ▶  $Q$  est l'espérance de la log-vraisemblance des données complètes  $\log(p(D, Z; \theta))$ , connaissant le modèle courant à l'étape  $t$ ,  $\theta(t)$
  - ▶ L'espérance est calculée par rapport à la distribution des variables cachées  $z$  connaissant le modèle courant  $\theta(t)$  :  $p(z|x; \theta(t))$
  - ▶  $Q(\theta, \theta(t)) = E_Z[\log p(D, Z; \theta); \theta(t)] = \sum_Z \log p(D, Z; \theta) p(Z|D; \theta(t))$
- ▶ Dans cette expression :
  - ▶  $D$  et  $\theta(t)$  sont des constantes
  - ▶  $z$  est une variable aléatoire de densité  $p(z|x; \theta(t))$ ,  $Z$  est l'ensemble des variables  $z$
  - ▶  $\theta$  représente les paramètres que l'on veut estimer

## Algorithme EM – Cas general

- ▶ On va montrer que  $Q(\theta, \theta(t))$  est une borne inférieure de  $L(\theta)$
- ▶ Pour cela on utilisera **l'inégalité de Jensen**
- ▶ Théorème
  - ▶ Soit  $f$  une fonction convexe définie sur  $R$  et  $x$  une variable aléatoire réelle, alors
  - ▶  $E[f(x)] \geq f(E[x])$
  - ▶ Si  $f$  est strictement convexe  $E[f(x)] = f(E[x])$  si et seulement si  $x$  est une constante
- ▶ Rq
  - ▶  $f$  est convexe si  $f''(x) \geq 0 \forall x$
  - ▶ Pour une fonction concave (e.g.  $\log()$ ), l'inégalité est vérifiée dans le sens opposé  $E[f(x)] \leq f(E[x])$

## Algorithme EM – Cas general

- ▶  $\log p(D; \theta) = \sum_{i=1}^N \log p(x^i; \theta)$
- ▶  $= \sum_{i=1}^N \log \sum_{z^i} p(x^i, z^i; \theta)$
- ▶  $= \sum_{i=1}^N \log \sum_{z^i} \frac{p(z^i|x^i; \theta(t))p(x^i, z^i; \theta)}{p(z^i|x^i; \theta(t))}$
- ▶  $= \sum_{i=1}^N \log E_{z^i \sim p(z|x; \theta(t))} \frac{p(x^i, z^i; \theta)}{p(z^i|x^i; \theta(t))}$
- ▶  $\geq \sum_{i=1}^N E_{z^i \sim p(z|x; \theta(t))} \log \frac{p(x^i, z^i; \theta)}{p(z^i|x^i; \theta(t))}$
- ▶  $= \sum_{i=1}^N \sum_{z^i} p(z^i|x^i; \theta(t)) \log \frac{p(x^i, z^i; \theta)}{p(z^i|x^i; \theta(t))}$
- ▶ On a utilisé la concavité de la fonction log pour l'inégalité
- ▶ Comme dans l'étape M, on maximise par rapport à  $\theta$ , on peut éliminer le dénominateur du log dans l'expression précédente, on a alors la fonction auxiliaire :
- ▶  $Q(\theta, \theta(t)) = \sum_{i=1}^N \sum_{z^i} p(z^i|x^i; \theta(t)) \log p(x^i, z^i; \theta)$
- ▶ En maximisant  $Q$ , on maximise une borne inférieure de la vraisemblance

## Algorithme EM – Cas general

- ▶ On peut montrer ensuite la convergence de l'algorithme par :
  - ▶  $Q(\theta(t+1), \theta(t)) \geq Q(\theta(t), \theta(t)) \Rightarrow p(D; \theta(t+1)) \geq p(D; \theta(t))$

## Algorithme EM - Cas général

Initialiser  $\theta(0)$

1. Etape E : *Espérance*

On calcule  $p(Z|D, \theta(t))$

On en déduit  $Q(\theta; \theta(t))$

L'espérance est calculée par rapport à la distribution de  $Z$  pour les paramètres courants  $\theta(t)$

2. Etape M : *Maximisation*

Etant donnée la distribution courante sur  $Z$ , trouver les paramètres qui maximisent  $Q$

$$\theta(t + 1) = \operatorname{argmax}_{\theta} E_{Z \sim p(Z|x;\theta(t))} [\log p(D, Z; \theta)]$$

► L'algorithme converge vers un maximum local de la fonction  $Q$  et de  $p(D; \theta)$



## Algorithme EM - Cas général

- ▶ Remarques
  - ▶ L'algorithme est utilisé pour
    - ▶ les algorithmes non supervisés, semi - supervisés
    - ▶ les données manquantes ou les composantes manquantes dans les données
    - ▶ les HMM ...

## Algorithme EM - Mélange de densités – cas gaussien

- ▶ On suppose que le modèle génératif des données est un mélange de densités gaussiennes

- ▶ On fixe a priori le nombre de composantes du mélange à  $k$
- ▶ Pour simplifier, on suppose que les données  $x$  sont unidimensionnelles

- ▶  $p(x) = \sum_{l=1}^k p(l)p(x|l), \quad p(x|l) = \frac{1}{(2\pi\sigma_l^2)^{\frac{1}{2}}} \exp\left(-\frac{(x-\mu_l)^2}{2\sigma_l^2}\right)$

- ▶ Paramètres

- ▶ Coefficients du mélange  $\phi_l = p(l)$ , moyennes  $\mu_l$  et écarts types  $\sigma_l$  :

$$\theta = \{p(l), \mu_l, \sigma_l; l = 1 \dots k\}$$

## Algorithme EM - Mélange de densités – cas gaussien

### ▶ Log-vraisemblance

$$\text{▶ } \log p(D; \theta) = \sum_{i=1}^N \log \sum_{l=1}^k p(l; \theta) p(x^i | l; \theta)$$

### ▶ Vraisemblance complète

$$\text{▶ } \log p(D, Z; \theta) = \sum_{i=1}^N \log(p(z^i; \theta) p(x^i | z^i; \theta)) = \sum_{i=1}^N \sum_{l=1}^k \delta_{lz^i} \log(p(l; \theta) p(x^i | l; \theta))$$

### ▶ Fonction auxiliaire

$$\text{▶ } Q(\theta, \theta(t)) = E_{z \sim p(z|x; \theta(t))} [\log p(D, Z; \theta)]$$

$$\text{▶ } = \sum_{z^1=1}^k \dots \sum_{z^N=1}^k \log p(D, Z; \theta) \prod_{i=1}^N p(z^i | x^i; \theta(t))$$

▶ quelques réécritures ... (cf slide suivant) conduisent à l'expression :

$$\text{▶ } = \sum_{i=1}^N \sum_{l=1}^k p(l | x^i; \theta(t)) \log(p(l; \theta) p(x^i | l; \theta))$$

▶ Qui sera utilisée dans la maximisation de l'étape M.

► Reécriture de la fonction auxiliaire

- $Q = Q(\theta, \theta(t)) = E_{\mathbf{z} \sim p(\mathbf{z}|\mathbf{x}; \theta(t))} [\log p(D, \mathbf{Z}; \theta)]$
- $Q = \sum_{z^1=1}^k \dots \sum_{z^N=1}^k \log p(D, \mathbf{Z}; \theta) \prod_{i=1}^N p(z^i | x^i; \theta(t))$
- $Q = \sum_{z^1=1}^k \dots \sum_{z^N=1}^k (\sum_{j=1}^N \sum_{l=1}^k \delta_{lz^j} \log(p(l; \theta)p(x^j | l; \theta))) \prod_{i=1}^N p(z^i | x^i; \theta(t))$
- $Q = \sum_{j=1}^N \sum_{l=1}^k \left[ \left( \sum_{z^1=1}^k \dots \sum_{z^N=1}^k \delta_{lz^j} \prod_{i=1}^N p(z^i | x^i; \theta(t)) \right) \log p(l; \theta)p(x^j | l; \theta) \right]$

- $A = \sum_{z^1=1}^k \dots \sum_{z^N=1}^k \delta_{lz^j} \prod_{i=1}^N p(z^i | x^i; \theta(t))$
- $A = \sum_{z^1=1}^k \dots \sum_{z^{j'-1}=1}^k \sum_{z^{j'+1}=1}^k \dots \sum_{z^N=1}^k (\sum_{z^{j'}=1}^k \delta_{lz^{j'}} \prod_{i=1}^N p(z^i | x^i; \theta(t)))$
- $A = \sum_{z^1=1}^k \dots \sum_{z^{j'-1}=1}^k \sum_{z^{j'+1}=1}^k \dots \sum_{z^N=1}^k (p(l | x^{j'}; \theta(t)) \prod_{i \neq j'} p(z^i | x^i; \theta(t)))$

►  $\sum_{z^1=1}^k \dots \sum_{z^{j'-1}=1}^k \sum_{z^{j'+1}=1}^k \dots \sum_{z^N=1}^k \prod_{i \neq j'} p(z^i | x^i; \theta(t)) = 1$

►  $Q = \sum_{j=1}^N \sum_{l=1}^k p(l | x^j; \theta(t)) \log p(l; \theta)p(x^j | l; \theta)$

## Algorithme EM - Mélange de densités – cas gaussien – Etapes E et M

### ▶ Etape E

$$\text{▶ } p(z^i = j | x^i; \theta(t)) = \frac{p(x^i | z^i = j; \theta(t)) p(z^i = j; \theta(t))}{\sum_{l=1}^k p(x^i | z^i = l; \theta(t)) p(z^i = l; \theta(t))} \quad \forall i, j$$

### ▶ Etape M

$$\text{▶ } \text{Min}_{\theta}(-Q(\theta, \theta(t))) \text{ sous la contrainte } \sum_{l=1}^k p(l; \theta) = 1$$

▶ Équivalent à

$$\text{▶ } \text{Min}_{\theta}(-Q(\theta, \theta(t)) + \lambda(\sum_{l=1}^k p(l; \theta) - 1))$$

▶ Avec  $\lambda$  le coefficient de Lagrange associé

### ▶ Remarque

▶ Dans l'étape E on n'a pas besoin de calculer explicitement  $Q(\theta, \theta(t))$ , il suffit de calculer les  $p(z^i = j | x^i; \theta(t))$  (cf dérivation du résultat en cours)

## Algorithme EM - Mélange de densités – cas gaussien – Reestimation dans l'étape M

- ▶ à l'étape M à l'itération  $t$ , on obtient les formules de reestimation suivantes

- ▶  $p(j) = \frac{1}{N} \sum_{i=1}^N p(z^i = j | x^i; \theta(t))$

- ▶  $\mu_j = \frac{\sum_{i=1}^N p(z^i = j | x^i; \theta(t)) x^i}{\sum_{i=1}^N p(z^i = j | x^i)}$

- ▶  $\sigma_j^2 = \frac{\sum_{i=1}^N p(z^i = j | x^i; \theta(t)) (x^i - \mu_j)^2}{\sum_{i=1}^N p(z^i = j | x^i)}$

## Algorithme EM - Mélange de densités – cas gaussien

### ▶ Etape E

- ▶ Si on connaît les lois de mélange (les  $(\mu, \sigma, \phi)$ ), on peut facilement calculer les  $p(z^i = j|x^i)$  par  $p(z^i = j|x^i) = \frac{p(x^i|z^i=j)p(z^i=j)}{p(x^i)}$

### ▶ Etape M

- ▶ Si on connaît les probabilités d'appartenance, on peut facilement estimer les paramètres des lois de mélange comme vu dans l'exemple
- ▶ Il suffit de partir d'une valeur initiale – pour les  $p(z^i = j|x^i)$  ou pour les paramètres des lois de mélange pour exécuter l'algorithme
  - ▶ Celui ci converge vers un maximum local de la vraisemblance

# Apprentissage non supervisé

Mélange de densités

Apprentissage par échantillonnage de Gibbs



# Les méthodes MCMC

## Markov Chain Monte Carlo

- ▶ Méthodes de calcul intensif basées sur la simulation pour
  - ▶ Echantillonnage de variables aléatoires
    - ▶  $\{x^t\}_{t=1..T}$  qui suivent une certaine distribution  $p(x)$
  - ▶ Calcul de l'espérance de fonctions suivant cette distribution
    - ▶  $E[f(x)]$  sera estimé par  $\frac{1}{T} \sum_{t=1}^T f(x^t)$
    - ▶ e.g. moyenne, marginales, ...
  - ▶ Maximisation de fonctions
    - ▶  $Argmax_x p(x)$

## Echantillonneur de Gibbs

- ▶ On veut estimer une densité  $p(\mathbf{x})$ ,  $\mathbf{x} \in R^n$ , avec  $\mathbf{x} = (x_1, \dots, x_n)$
- ▶ Hypothèse
  - ▶ On connaît les lois conditionnelles
    - ▶  $p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$  qui sera notée  $p(x_i | x_{-i})$
- ▶ Algorithme
  - ▶ Initialiser  $\{x_i^0, i = 1, \dots, n\}$
  - ▶ Itérer jusqu'à convergence (variable d'itération notée  $t$ )
  - ▶ échantillonner
    - ▶  $x_1^{t+1} \sim p(x_1 | x_{-1}^t) = p(x_1 | x_2^t, \dots, x_n^t)$
    - ▶ .....
    - ▶  $x_n^{t+1} \sim p(x_n | x_{-n}^t) = p(x_n | x_1^{t+1}, \dots, x_{n-1}^t)$
  - ▶ Jusqu'à ce que la distribution jointe  $p(x_1, \dots, x_n)$  ne change plus

## ▶ Propriétés

- ▶ Sous certaines conditions de régularité, la procédure converge vers la distribution cible  $p(\mathbf{x})$ 
  - ▶ Les échantillons résultants sont des échantillons de la loi jointe  $p(\mathbf{x})$
- ▶ On n'a pas besoin de connaître la forme analytique des  $p(x_i|x_{-i})$  mais uniquement de pouvoir échantillonner à partir de ces distributions
  - ▶ Mais la forme analytique permet d'avoir de meilleurs estimés
- ▶ Avant de retenir les points échantillons, on autorise souvent une période de “burn-in” pendant laquelle on fait simplement tourner l'algorithme “à vide”
- ▶ Gibbs facile à implémenter, adapté aux modèles hiérarchiques (cf LDA)

# Cas du mélange de deux lois gaussiennes

- ▶ Modèle :  $p(x) = \sum_{l=1}^2 p_l p(x|l)$ 
  - ▶ On considère un mélange de deux gaussiennes à une dimension
- ▶ On va considérer un modèle augmenté en ajoutant une variable cachée  $z \in \{0,1\}$ 
  - ▶ Les données complètes sont les  $(x^i, z^i)$
- ▶ Les paramètres à estimer sont
  - ▶ Les paramètres des gaussiennes :  $\theta = \{p_l, \mu_l, \sigma_l; l = 1,2\}$
  - ▶ La variable cachée  $z$  (qui est considérée comme un paramètre supplémentaire à estimer dans la procédure de gibbs)
- ▶ On va utiliser Gibbs en échantillonnant sur les densités conditionnelles des paramètres  $p(\theta, z|x)$ 
  - ▶ Pour simplifier on suppose dans l'exemple que les proportions  $p_l$  et les variances  $\sigma_l$  sont fixées, on estime juste les moyennes  $\mu_1, \mu_2$
- ▶ Pour cela, on va échantillonner suivant la distribution jointe  $p((\mu_1, \mu_2), z|x)$

## Echantillonneur de Gibbs pour le modèle de mélange de deux gaussiennes

- ▶ Hyp:  $p_i$  et  $\sigma_i, i = 1, 2$  sont supposés connus
  - ▶ Choisir des valeurs initiales  $\mu_1(0), \mu_2(0)$
  - ▶ Répéter ( $t = 1, 2, \dots$ ) jusqu'à convergence (i.e.  $\mu_j(t), j = 1, 2$  ne changent plus)
    - ▶ Pour  $i = 1 \dots N$ 
      - Générer  $z^i(t) \in \{0, 1\}$  selon la distribution
        - $p(z^i = 1 | x^i) = \frac{p_1 p(x^i | z^i=1; \mu_1(t-1), \sigma_1)}{p_1 p(x^i | z^i=1; \mu_1(t-1), \sigma_1) + p_2 p(x^i | z^i=0; \mu_2(t-1), \sigma_2)}$
        - Rq: c'est une loi de Bernouilli de paramètre le membre droit de l'équation
      - ▶ Calculer
        - $\hat{\mu}_1(t) = \frac{\sum_{i=1}^N z^i(t) x^i}{\sum_{i=1}^N z^i(t)}$
        - $\hat{\mu}_2(t) = \frac{\sum_{i=1}^N (1 - z^i(t)) x^i}{\sum_{i=1}^N (1 - z^i(t))}$
- Et générer  $\mu_j(t) \sim N(\hat{\mu}_j, \sigma_j), j = 1, 2$

## Lien avec l'algorithme EM

- ▶ Les étapes pour cet exemple sont les mêmes qu'avec EM
- ▶ Différence
  - ▶ Au lieu de maximiser la vraisemblance, aux étapes 1 et 2, on échantillonne
    - ▶ Etape 1 : on simule les variables cachées  $z$  à partir de la distribution conditionnelle  $p(z|\theta, x)$  au lieu de calculer  $E(z|\theta, x)$  dans EM
    - ▶ Etape 2 : on simule à partir de  $p(\mu_1, \mu_2 | z, x)$  au lieu de calculer le max. vraisemblance  $p(\mu_1, \mu_2, z|x)$  dans EM

## ▶ Remarque 1

- ▶ Pour simplifier on a supposé que les proportions  $p_i$  et les variances  $\sigma_i$ ,  $i = 1, 2$  étaient fixes.
- ▶ Dans le cas général où elles ne sont pas connues, on va
  - ▶ 1. fixer des valeurs initiales de l'ensemble des paramètres  $\theta = \{p_i, \mu_i, \sigma_i; i = 1, 2\}$  à partir de distribution a priori pour ces paramètres
  - ▶ 2. échantillonner alternativement pour chaque paramètre selon sa distribution a posteriori, conditionnellement aux autres paramètres
    - Par exemple on échantillonnera successivement pour  $(p_1, p_2)$ ,  $(\mu_1, \mu_2)$ , puis  $(\sigma_1, \sigma_2)$

## ▶ Remarque 2

- ▶ Il existe d'autres méthodes pour estimer des densités de façon approchée quand il n'existe pas de forme analytique conduisant à des calculs, e.g. méthodes variationnelles

# Apprentissage non supervisé

Spectral Clustering



## Spectral Clustering (after Von Luxburg 2007)

### ▶ Intuition

- ▶  $\mathbf{x}_1, \dots, \mathbf{x}_n$  data points,  $w_{ij}$  similarity between  $\mathbf{x}_i$  and  $\mathbf{x}_j$
- ▶  $G = (V, E)$  undirected graph
  - ▶  $V = \{v_1, \dots, v_n\}$ , vertex  $v_i$  corresponds to data point  $\mathbf{x}_i$
  - ▶ Edges are weighted,  $W = (w_{ij})_{i,j=1\dots n}$ ,  $w_{ij} \geq 0$  is the weight matrix
  - ▶  $D$  : diagonal matrix with  $d_i = \sum_{j=1}^n w_{ij}$
- ▶ Clustering amounts at finding a graph partition such that
  - ▶ Edges between clusters have low weights
  - ▶ Edges among points inside a cluster have high values

## Spectral Clustering

- ▶ **Building similarity graphs from data points**
  - ▶ Different ways to build a similarity graph
  - ▶ Build a locally connected graphs: k-nearest neighbor graphs
    - ▶ Two vertices are connected if one of them is among the k-nearest neighbor of the other
    - ▶ Or two vertices are connected if both are in the k-neighborhood of the other
    - ▶ Edges are then weighted using the similarity of the vertices
  - ▶ Build a fully connected graphs
    - ▶  $w_{ij} = \exp(-\|x_i - x_j\|^2 / 2\sigma^2)$

# Spectral Clustering

## ▶ Graph Laplacians

### ▶ Unnormalized graph Laplacian

- ▶  $L = D - W$

### ▶ Normalized graph Laplacians

- ▶  $L_{sym} = D^{-\frac{1}{2}}LD^{-\frac{1}{2}} = I - D^{-\frac{1}{2}}WD^{-\frac{1}{2}}$  symmetric

- ▶  $L_{rw} = D^{-1}L = I - D^{-1}W$  interpretation : random walk on the graph

## Spectral Clustering

### ▶ Properties of the unnormalized graph Laplacian $L$

#### ▶ Proposition 1- $L$ satisfies:

- ▶  $\forall y \in R^n, \quad y^T L y = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (y_i - y_j)^2$
- ▶  $L$  is symmetric, positive semi-definite
- ▶ The smallest eigenvalue of  $L$  is 0, the corresponding eigenvector is  $\mathbf{1}$  (vector with  $n$  1s)
- ▶  $L$  has  $n$  non negative eigenvalues  $0 = \lambda_1 \leq \dots \leq \lambda_n$

#### ▶ Proof

- ▶  $y^T L y = y^T D y - y^T W y = \sum_{i=1}^n d_i y_i^2 - \sum_{i,j=1}^n y_i y_j w_{ij} = \frac{1}{2} (\sum_{i=1}^n d_i y_i^2 - 2 \sum_{i,j=1}^n y_i y_j w_{ij} + \sum_{j=1}^n d_j y_j^2) = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (y_i - y_j)^2$
- ▶ Symmetry follows from the symmetry of  $L$  and  $W$ , positive semi definiteness comes from  $y^T L y \geq 0$
- ▶ Obvious: a symmetric, positive semi-definite matrix has all its eigenvalues  $\geq 0$
- ▶ Direct consequence of properties 1 and 3

## Spectral Clustering

- ▶ **Proposition 2 (number of connected components and spectrum of  $L$ )**
  - ▶ Let  $G$  be an undirected graph with non negative weights. The multiplicity  $k$  of the eigenvalue  $0$  of  $L$  equals the number of connected components in the graph. The eigenspace of eigenvalue  $0$  is spanned by the indicator vectors of these components.
  - ▶ This result provides intuition on spectral clustering algorithms
  - ▶ **Connected graph**
    - ▶ Let us consider first  $k = 1$ . If  $y$  is an eigenvector with eigenvalue  $0$  then  $0 = y^T L y = \sum_{i=1}^n w_{ij} (y_i - y_j)^2$ . Since  $w_{ij} \geq 0$  this is only possible if  $y_i = y_j \forall i, j$ . In a graph with 1 connected component,  $\mathbf{1}$  the constant vector filled with 1s is the eigenvector with eigenvalue  $0$ . It is also the indicator vector of the connected component
  - ▶ **Multiple components**
    - ▶ If there are  $k$  connected components. The  $W$  matrix and hence  $L$  can be arranged in a block diagonal form. The above result holds for each connected component  $A_i$ . The corresponding eigenvectors are the  $\mathbf{1}_{A_i}$  with 1s corresponding to the  $i^{th}$  block and 0s everywhere else.
    - ▶ If one can identify these eigenvectors, one will identify the clusters
      - This is only intuition, the situation is usually more complex (see Luxburg 2007)

## Spectral Clustering

### ► Properties of the normalized graph Laplacians

- $\forall y \in R^n, \quad y^T L_{sym} y = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \left( \frac{y_i}{\sqrt{d_i}} - \frac{y_j}{\sqrt{d_j}} \right)^2$
- $L_{sym}$  and  $L_{rw}$  are positive semi-definite and have  $n$  non negative eigenvalues  
 $0 = \lambda_1 \leq \dots \leq \lambda_n$
- $\lambda$  is an eigenvalue of  $L_{rw}$  with eigenvector  $u$  iff  $\lambda$  is an eigenvalue of  $L_{sym}$  with eigenvector  $D^{1/2}u$

## Unnormalized spectral clustering algorithm

### ▶ Idea

- ▶ Project data points  $\mathbf{x}_i \in R^m, i = 1 \dots n$ , in a smaller dimensional space, say of dimension  $k$  where clustering is performed

### ▶ Input: $n$ points $\mathbf{x}_1, \dots, \mathbf{x}_n$ , similarity matrix $S$

### ▶ Output: clusters

- ▶ Construct similarity graph and corresponding weight matrix  $W$
- ▶ Compute unnormalized Laplacian  $L$
- ▶ Compute  $k$  first eigenvectors of  $L$  (corresponding to smallest eigenvalues):  $\mathbf{u}_1, \dots, \mathbf{u}_k$
- ▶  $U$ :  $n \times k$  matrix with columns  $\mathbf{u}_1, \dots, \mathbf{u}_k$
- ▶ For  $i = 1 \dots n, \mathbf{y}_i \in R^k$   $i$ -th row of  $U$
- ▶ Cluster the  $\mathbf{y}_i, i = 1 \dots n$  with  $k$ -means algorithm into clusters  $C_1, \dots, C_k$
- ▶  $k$  clusters in the initial space:  $C'_1, \dots, C'_k / C'_i = \{\mathbf{x}_j / \mathbf{y}_j \in C_i\}$

### ▶ Note: Similar algorithms with normalized Laplacians

# Apprentissage non supervisé

Non Negative Matrix Factorization



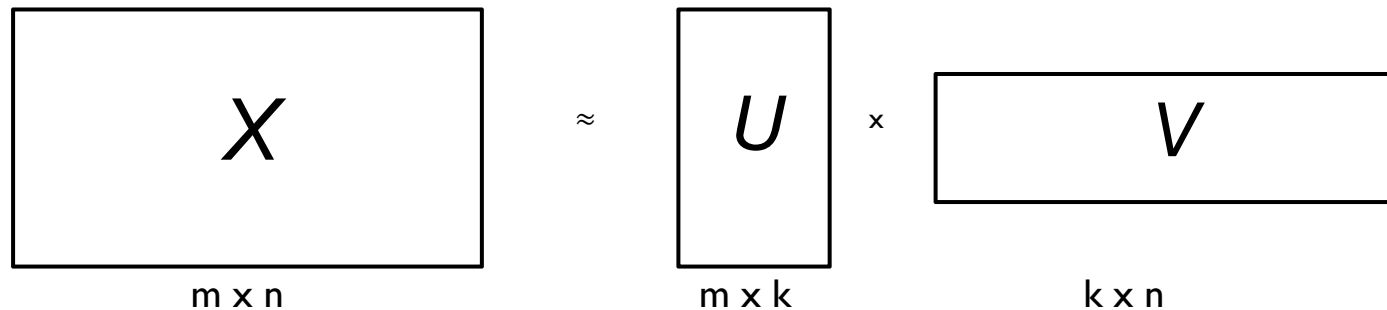
## Matrix Factorization

### ▶ Idea

- ▶ Project data vectors in a latent space of dimension  $k < m$  size of the original space
- ▶ Axis in this latent space represent a new basis for data representation
- ▶ Each original data vector will be approximated as a linear combination of  $k$  basis vectors in this new space
- ▶ Data are assigned to the nearest axis
- ▶ This provide a clustering of the data

## Matrix factorization

- ▶  $X = \{x_{.1}, \dots, x_{.n}\}, x_i \in R^m$
- ▶  $X$   $m \times n$  matrix with columns the  $x_{.i}$  s
- ▶ Low rank approximation of  $X$ 
  - ▶ Find factors  $U, V, / X \approx UV$
  - ▶ With  $U$  an  $m \times k$  matrix,  $V$  a  $k \times n$  matrix,  $k < m, n$



- ▶ Many different decompositions
  - ▶ e.g. Singular Value Decomposition, Non Negative Matrix Factorization, Tri factorization, etc

## ▶ Applications

- ▶ Recommendation (User x Item matrix)
  - ▶ Matrix completion

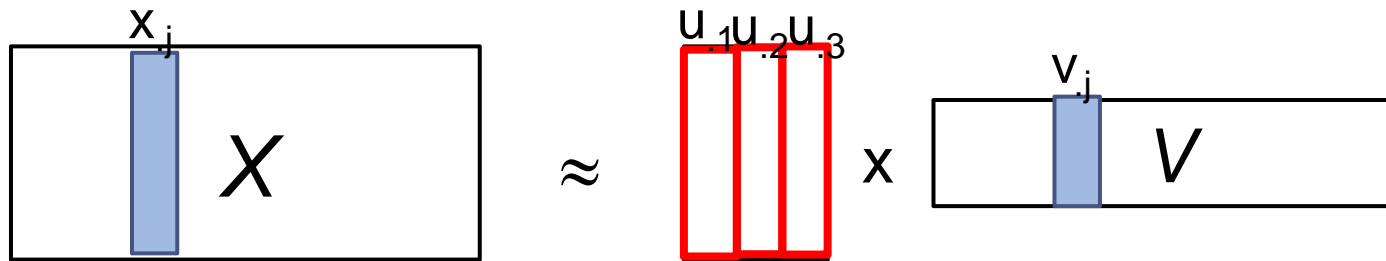
	Star wars	Terminator	Titanic	Scarface
Bob	5	4	3	?
Alice	4	?	5	3
John	?	5	3	4

- ▶ Link prediction (Adjacency matrix)
- ▶ ...

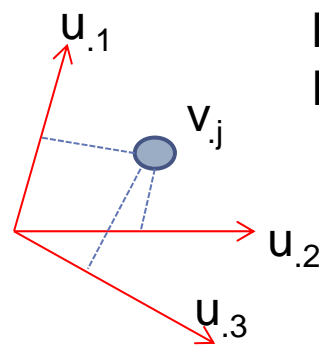
▶  $X \approx UV$

▶  $x_j = \sum_{i=1}^k u_i v_{ij}$

▶ Columns of  $U$ ,  $u_j$  are basis vectors, the  $v_{ij}$  are the coefficient of  $x_j$  in this basis



Original data

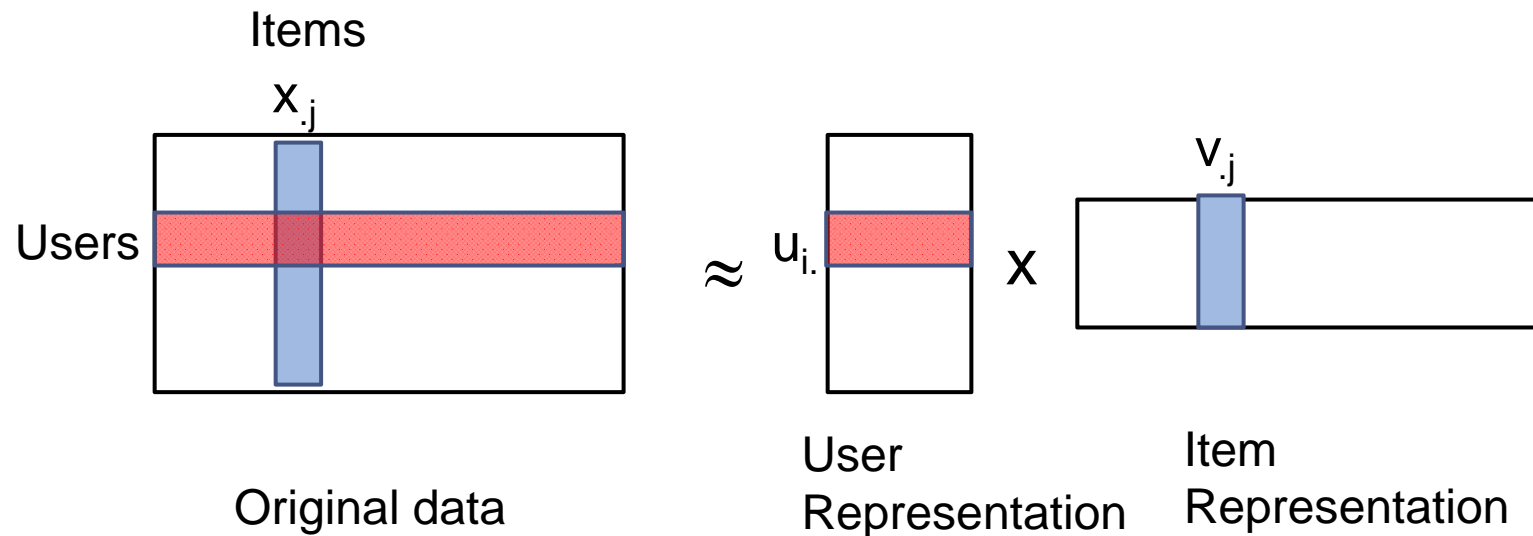


Basis vectors  
Dictionary Representation



► Interpretation

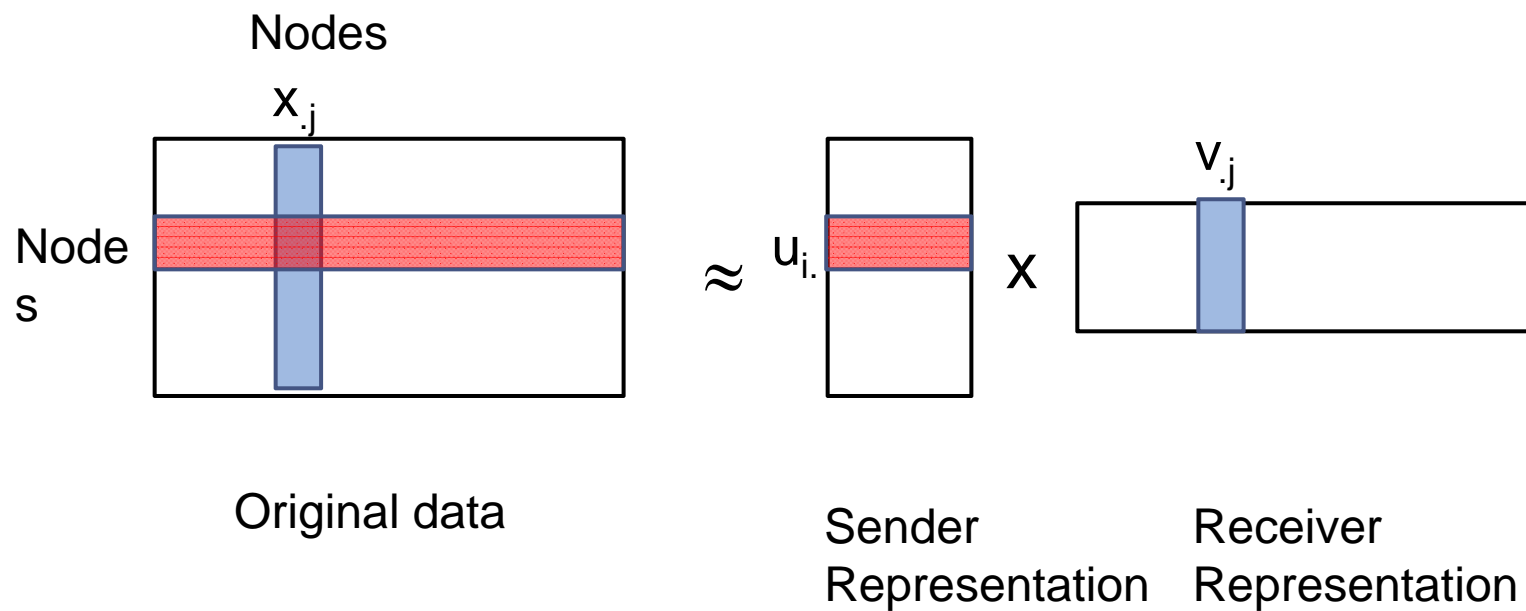
- If  $X$  is a User  $\times$  Item matrix



- Users and items are represented in a common representation space of size  $k$
- Their interaction is measured by a dot product in this space

► Interpretation

- If  $X$  is a directed graph adjacency or weight matrix



▶ Loss function - example

▶ Minimize  $C = \|X - UV\|^2 + c(U, V)$

▶ constraints on  $U, V$ :  $c(U, V)$  e.g.

- ▶ Positivity (NMF)
- ▶ Sparsity of representations, e.g.  $\|V\|_1$ , group sparsity, ...
- ▶ Overcomplete dictionary  $U$ :  $k > m$
- ▶ Symmetry, e.g. trifactorisation
- ▶ Bias on  $U$  and  $V$ 
  - e.g. biases of users (low scores) or item (popularity) for recommendation
- ▶ Multiple graphs
- ▶ Any a priori knowledge on  $U$  and  $V$

## Non Negative Matrix Factorization

- ▶ Loss function

- ▶ Solve

$$\text{Min}_{\{U,V\}} \|X - UV\|^2$$

Under constraints  $U, V \geq 0$

i.e. the factors are constrained to be non negative

- ▶ Convex loss function in  $U$  and in  $V$ , but not in both  $U$  and  $V$



## ▶ Algorithm

- ▶ Constrained optimization problem
- ▶ Can be solved by a Lagrangian formulation
  - ▶ Iterative multiplicative algorithm (Xu et al. 2003)
    - U,V initialized at random values
    - Iterate until convergence
      - $u_{ij} \leftarrow u_{ij} \frac{(XV^T)_{ij}}{(UVV^T)_{ij}}$
      - $v_{ij} \leftarrow v_{ij} \frac{(X^TU)_{ij}}{(V^TU^TU)_{ij}}$
- ▶ Or by projected gradient formulations
- ▶ The solution  $U, V$  is not unique, if  $U, V$  is solution, then  $UD, D^{-1}V$  for  $D$  diagonal positive is also solution

## ▶ Using NMF for Clustering

- ▶ Normalize  $U$  as a column stochastic matrix (each column vector is of norm 1)

- ▶  $u_{ij} \leftarrow \frac{u_{ij}}{\sqrt{\sum_i u_{ij}^2}}$

- ▶  $v_{ij} \leftarrow v_{ij} \sqrt{\sum_i u_{ij}^2}$

- ▶  $U$  columns are cluster centroids,  $V$  columns are cluster membership indicators.
- ▶ Under the constraint “ $U$  normalized” the solution  $U, V$  is unique
- ▶ Associate  $x^i$  to cluster  $j$  if  $j = \operatorname{argmax}_j(v_{ji})$

## ▶ Note

- ▶ Many different versions and extensions of NMF
- ▶ Different loss functions
  - ▶ e.g. different constraints on the decomposition
- ▶ Different algorithms

## ▶ Applications

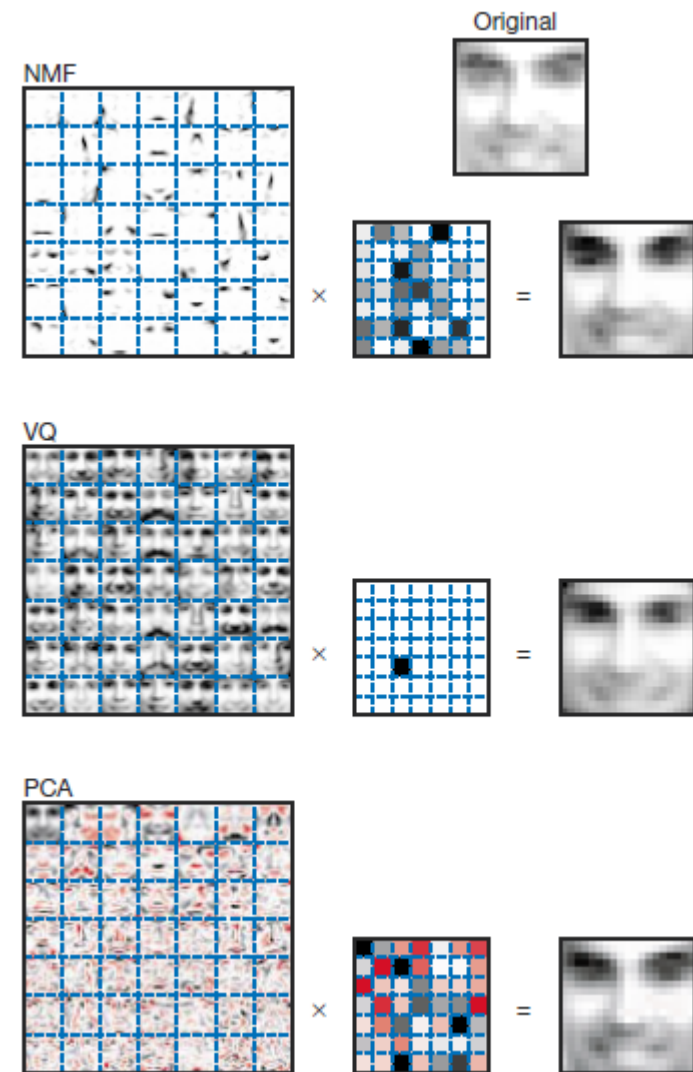
- ▶ Clustering
- ▶ Recommendation
- ▶ Link prediction
- ▶ Etc

## ▶ Specific forms of NMF can be shown equivalent to

- ▶ PLSA
- ▶ Spectral clustering

## Illustration (Lee & Seung 1999)

- ▶ Basis images for
- ▶ NMF
- ▶ Vector Quantization
- ▶ Principal Component Analysis



## ▶ Citations

- ▶ Ulrike Luxburg. 2007. A tutorial on spectral clustering. *Statistics and Computing* 17, 4 (December 2007), 395-416.
- ▶ Daniel D. Lee and [H. Sebastian Seung](#) (1999). "Learning the parts of objects by non-negative matrix factorization". *Nature* 401 (6755): 788–791.
- ▶ W. Xu, X. Liu, and Y. Gong. (2003) Document clustering based on non-negative matrix factorization. *SIGIR '03*, page 267
- ▶ Yang and E. Oja. Linear and nonlinear projective nonnegative matrix factorization. *IEEE Trans. on Neural Networks*, 21:734–749, 2010.