

Pattern Analysis 02.05.13

Mean shift Algorithm is a mode-seeking algorithm

=> Related application: clustering -

PDF



K-means:

- obtain (guess) initial distribution of cluster centers
- for each data point: identify the closest cluster center
- each cluster center is replaced by the coordinate-wise average of all data points, that are closest to it

=> REPEAT (until convergence)

[Convergence is guaranteed  
at the local minimum  
depends on the initialization]

GREEDY ALGORITHM / SEARCH:

is a strictly local optimization

Advantage: speed

Disadvantage: No brain

\* euclidean

Book: elements of statistical learning  
anhang 2 ver 10 . pdf

(2)  
Pattern Analysis 03.05.19

Two very common (dis-)similarity criteria for clusterings is the "within" or "intra" cluster distance  $w(c)$   
 $\Rightarrow$  Eq. 14.28) and the "between" or "inter" cluster distance  
 $B(c)$  ( $\Rightarrow$  Eq. 14.29)

$\rightarrow$  the k-Means algorithm minimizes the within-cluster distances  
 (greedily) ( $\rightarrow$  Eq. 14.31 - 14.33)  
 ↑  
 local minimum

Question: How can we determine a reasonable parameter  $k$ ?

$\Rightarrow$  Low-dimensional dataset: Look at the data

$\Rightarrow$  High-dimensional dataset: Need a mathematical criterion  
 ↓

- k-means minimizes  $w(c)$

-  $w(c)$  decreases for increasing  $k$ , for  $k=N$ :  $w(c)=0$   
 "every point is in its own cluster"

- one trick to determine  $K$  is Tibshirani's "gap statistics":

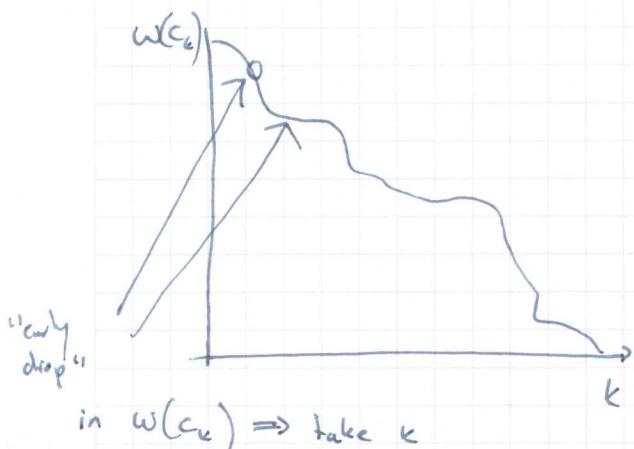
stop at the  $K$  where  $G(k) \geq G(k+1) - s'_{k+1}$

where  $G(k) = \log(w(c_k)) - \log(w(c_1))$   
 ↑  
 negative number

$s'_{k+1}$  is the standard deviation of the outcomes for  $k+1$  clusters

## Pattern Analysis (3)

03.05.14

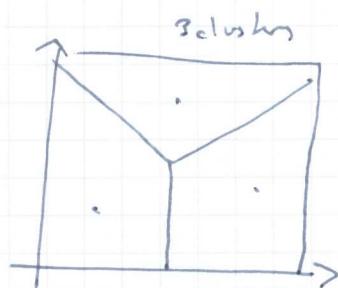
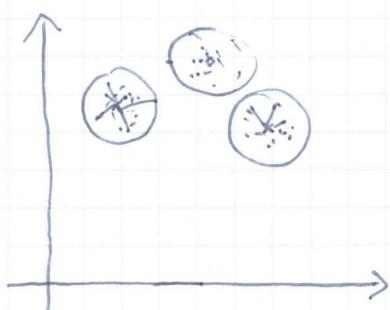


Book reference  
~ page 520

59:13

- A second trick is to create a reference curve from the uniform distribution : consider  $w(c_k^{\text{uniform}})$

$\Rightarrow$  compute the ratio of  $w(c_k)$  over  $w(c_k^{\text{uniform}})$   
and pick the minimum (or an early minimum).



good match  $\rightarrow$  relatively low intra cluster distance

Pattern Analysis

(4) / (1)

03.05.19 / 08.05.19

Example for agglomerative clustering: Felzenszwalb &amp; Huttenlocher

"Efficient Graph-based Image Segmentation"

08.05.19

Clustering[Paper] Ng, Jordan, Weiss : "on spectral clustering"Algorithm: Given points  $S = \{\vec{s}_1, \dots, \vec{s}_n\}, \vec{s}_i \in \mathbb{R}^d$ 1.) Compute affinities  $A \in \mathbb{R}^{n \times n}, A_{ij} = e^{-\frac{\|s_i - s_j\|_2^2}{2\sigma^2}}, A_{ii} = \emptyset$   
(diagonal is zero)2.) Let  $D = \text{diag}(\sum_j A_{1,j}, \sum_j A_{2,j}, \dots, \sum_j A_{n,j})$  of the row affinities,and set  $L = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$ 

↑  
Graph Laplacian

3.) Compute eigendecomposition of  $L$ Let  $X = (\vec{x}_1 \vec{x}_2 \vec{x}_3 \dots \vec{x}_k), X \in \mathbb{R}^{n \times k}$  denote a matrix consistingof the  $k$  largest eigenvectors  $\vec{x}_1, \dots, \vec{x}_k$  of  $L$ 

↖  
"eigenvectors associated with the largest  
eigenvalues"

4.) Normalize  $X$ :  $y_{ij} = \frac{x_{ij}}{\sqrt{(\sum_i x_{ij}^2)^{1/2}}}$ 5.) Perform k-means clustering with  $k$  centers on the points. Then are the rows of  $Y$ 6.)  $\Rightarrow$  The points in the original space are assigned to the  $k$  clusters  
of their transformed version in  $Y$ .

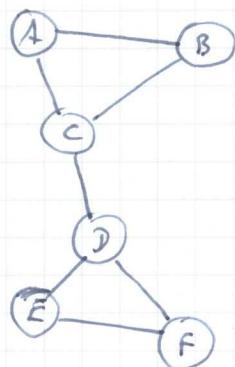
## Pattern Analysis ②

08.05.19

Quick, informal note on the Graph Laplacian:

A tool that links graph theory, linear algebra and probability theory.

Let's look at a graph:



Popular graph algorithms:

Maximum Flow / shortest Path / ... / between two nodes?

The set of mathematical tools, that can be used on a graph is somewhat limited: what exactly should we calculate on a list of nodes and edges?

⇒ Spectral Graph Theory opens a door to linear algebra, by treating a graph's adjacency matrix as an actual matrix in an algebraic sense ⇒ we can - for example - compute the eigenvalues to find out invariants about the underlying graph

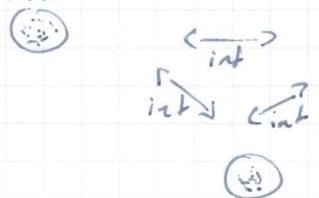
⇒ A useful representation of the adjacency matrix is the Graph Laplacian L

## Pattern Analysis

(3)

08.05.19

1-16



Remarks on the "ideal case" (sec. 3.1):

- The set of eigenvalues of a block matrix  $\hat{L} = \begin{pmatrix} L^{(11)} & & \\ & L^{(22)} & \\ & & L^{(33)} \end{pmatrix}$

is identical to the set of  $\{ \text{eigenvalues}(L^{(11)}) , \text{eigenvalues}(L^{(22)}) , \text{eigenvalues}(L^{(33)}) \}$

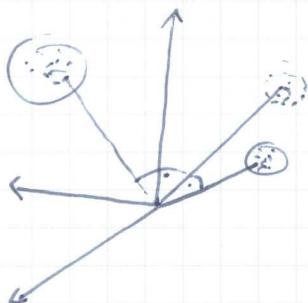
(can be treated independently)

The largest eigenvalue of a graph laplacian is 1, the second largest is lower than 1

$\Rightarrow \hat{L}$  has exactly 3 eigenvalues that are 1. These are the 3 largest ones, that we are looking for.

The associated 3 eigenvectors are orthogonal and thus span a 3-dimensional sub-space.

[Note], however, that the exact choice of eigenvectors is not unique



## Pattern Analysis (4)

08.05.19

The points are clustering around the rotation axes, that span the sub-space. These axes are spanned by  $90^\circ$  angles between them. Therefore k-Means can easily cluster the points in that subspace.

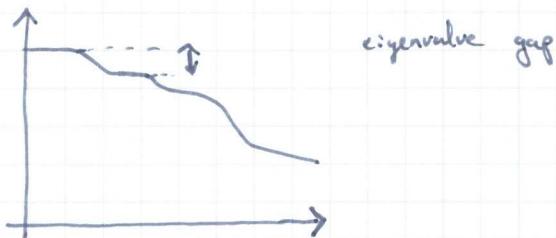
Sec 3.2, clusters that lie close to one another

The distribution of eigenvalues indicates, how well the data structures into k clusters : consider the eigenvalue gap

$$\delta = |\lambda_k - \lambda_{k+1}|$$

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \lambda_4 \geq \lambda_5 \geq \dots \geq$$

1 1 1 1-5



Cheeger constant:

$$h(S_i) = \min_I \frac{\sum_{j \in I} A_j}{\min \left\{ \sum_{j \in I} d_{i,j}^{(i)}, \sum_{k \notin I} d_k^{(i)} \right\}}$$

"between cluster  
distance"

Pattern Analysis

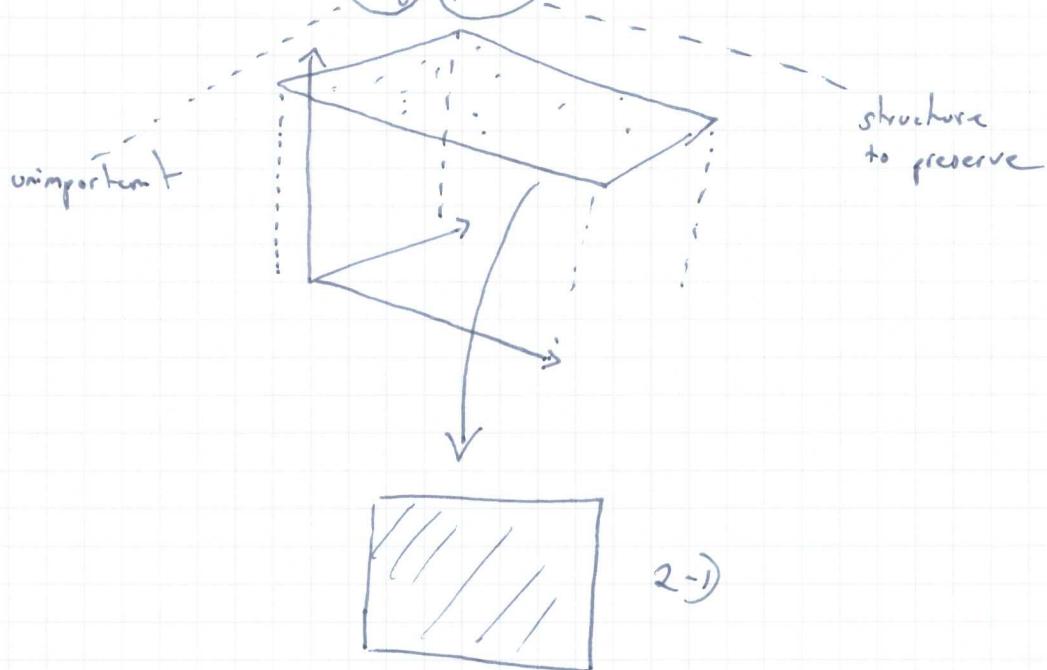
(1)

17.05.19

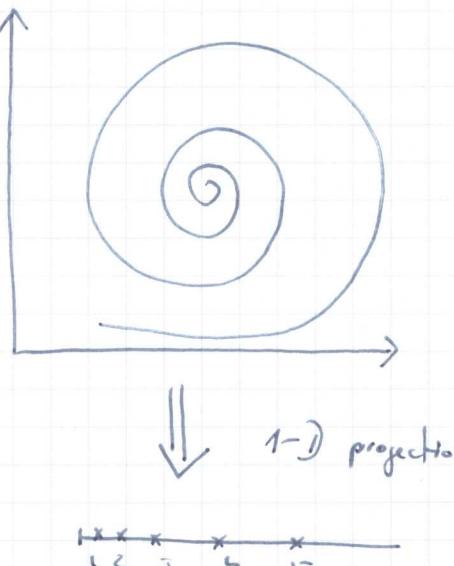
Manifold Learning -

Purpose / goal : Reduce Dimensionality of the data, while preserving its structure.

Example: consider a noisy plane in 3-D :



2nd example :



"swiss roll"  
↑  
this is the structure!

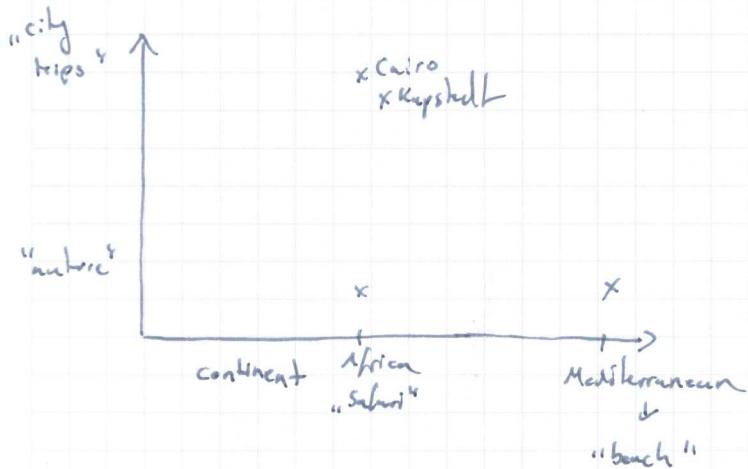
"unrolling"

## Pattern Analysis ②

17.05.19

In many real world tasks the structure to preserve and the low-dim manifold is not super-sharply defined.

Consider for example the manifold of tourist pictures



„Semantic Image Classification“

From millions of pixels in  
millions of images to very  
few ( $< 10$ ) dimensions  
 $\Rightarrow$  rather extreme task

In this class I would like to see manifold learning as the task of finding lower-dimensional representation of structure, essentially by ignoring the noise.

Why should we reduce the dimensionality?

$\Rightarrow$  To avoid the “curse of dimensionality”

Curse of Dimensionality - : Distance metrics tend to “wash out” in higher dimensional spaces, i.e. to lose their discriminative power.

To see this, consider  $N$  features  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ , where

$$0 \leq x_{i,k} \leq 1$$

$\vec{x}_{ik}$  kth element of the i-th vector

## Pattern Analysis (3)

17.05.19

To capture a fraction  $r$  of uniformly distributed features, we need to consider  $r \cdot V$  of the Volume  $V$  of the feature space.

This corresponds to a hypercube with  $d$  dimensions and edge length  $e_d(r) = r^{\frac{1}{d}} = \sqrt[d]{r}$

For example, to find 1% of the features in a 10-dim space,  $e_{10}(0,01) = (0,01)^{\frac{1}{10}} = 0,63$

to find 10% in 10-dim space:

$$e_{10}(0,1) = 0,1^{\frac{1}{10}} \approx 0,8$$

The median distance of the nearest neighbor to the origin of a  $d$ -dim space with  $N$  samples is

$$d(d, N) = \left(1 - \left(\frac{1}{2}\right)^{\frac{1}{N}}\right)^{\frac{1}{d}}$$

$\Rightarrow d(10, 5000) = 0,52$  most points are close to the boundary in some dimension(s)

PCA (Principal Component Analysis)

Find a linear mapping  $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}, d' \ll d$

that maximizes the variance in each dimension

## Pattern Analysis ④

17.05.19

objective Function  $J = \sum_{i,j=1}^N (\Phi \cdot \vec{x}_i - \Phi \vec{x}_j)^T$

$x \in \mathbb{R}^d$

$\cdot (\Phi \cdot \vec{x}_i - \Phi \vec{x}_j) + \\ + \lambda (\Phi^T \Phi - 1)$

$\Rightarrow$  PCA requires an Eigenvalue decomposition of the covariance matrix, PCA generates a low-dimensional orthogonal basis.

Multi-dimensional scaling (MDS)

Goal: To compute a low-dimensional representation of data points where we only know the distances (or dissimilarities) between them.

	Nürnberg	Beijing	Berlin	Mumbai
Nürnberg	-	9000	482,3	832,5
Beijing		-	8500	4000
Berlin			-	7900
Mumbai				-

$\Rightarrow$  can we reconstruct a map from the distances?

( $\rightarrow$  note, that there are close links to kernel PCA!!)

$\Rightarrow$  see ESSL.pdf, search for "MDSR"

Exercise 18.15

## Pattern Analysis (5)

17.05.19

Let  $X = (\vec{x}_1, \vec{x}_2, \vec{x}_3) \in \mathbb{R}^{d \times N}$  denote the feature vectors (as usual)

Set  $B = X^T X$

Let furthermore denote  $D^2 = [d_{ij}^2]_{i,j \in \{1, \dots, N\}}$

where  $d_{ij}^2 = (\vec{x}_i - \vec{x}_j)^T (\vec{x}_i - \vec{x}_j)$  (Euclidean Distance)

(could be substituted  
for another  
metric)

Task: given  $D^2$ , compute  $X$

$$d_{ij}^2 = (\vec{x}_i - \vec{x}_j)^T (\vec{x}_i - \vec{x}_j) = \vec{x}_i^T \vec{x}_i + \vec{x}_j^T \vec{x}_j - 2 \vec{x}_i^T \vec{x}_j$$

Assume that  $\vec{x}_1, \dots, \vec{x}_N$  are zero-mean, i.e.  $\sum_{i=1}^N \vec{x}_i = \emptyset$

The squared distance matrix is  $D^2 = \text{diag}(X^T X) \cdot \vec{1}^T + \vec{1} \cdot \text{diag}(X^T X) + 2 X^T X$   
where  $\vec{1} = (1, \dots, 1)^T \in \mathbb{R}^N$

"Magic matrix" / centering matrix  $C = (I - \frac{1}{N} \vec{1} \cdot \vec{1}^T)$

$\Rightarrow$  Multiplying  $D^2$  by the centering matrix from left & right

and weighing the result by  $-\frac{1}{2}$  yields:

$$-\frac{1}{2} C D^2 C = -\frac{1}{2} \left( I - \frac{1}{N} \vec{1} \cdot \vec{1}^T \right) \cdot \left( \text{diag}(X^T X) \cdot \vec{1}^T + \vec{1} \cdot \text{diag}(X^T X) + 2 X^T X \right) \cdot \left( I - \frac{1}{N} \vec{1} \cdot \vec{1}^T \right)$$

## Pattern Analysis (6)

17.05.14

$$\text{Term (1)} : \left[ \left( I - \frac{1}{N} \vec{1} \vec{1}^T \right) \cdot \text{diag}(x^T x) \right] \cdot \underbrace{\left[ \vec{1}^T \cdot \left( I - \frac{1}{N} \cdot \vec{1} \vec{1}^T \right) \right]}_{= \emptyset} = \emptyset$$

$$\vec{1}^T \cdot I - \frac{1}{N} \underbrace{\vec{1}^T \vec{1}}_N \cdot \vec{1}^T = \vec{1}^T - \vec{1}^T = 0$$

$$\text{Term (2)} : \left[ \text{analogously} \right] = \emptyset$$

$$\text{Term (3)} : \left[ -\frac{1}{2} \left( I - \frac{1}{N} \vec{1} \vec{1}^T \right) \cdot (-2x^T x) \cdot \left( I - \frac{1}{N} \vec{1} \vec{1}^T \right) \right] =$$

$$= I \cdot x^T - \frac{1}{N} \underbrace{\left( \vec{1} \vec{1}^T \right) x^T}_{\emptyset \text{ because of zero-mean of feature vectors}} \cdot \left( x - \frac{1}{N} x \cdot \vec{1} \vec{1}^T \right) \Rightarrow \boxed{x^T x = B}$$

!

$\Rightarrow$  Factorize  $B$  to obtain  $X$  via SVD

$$B = U^T \Sigma V$$

$$B \text{ is symmetric} \dots \Rightarrow \underbrace{U \Sigma^{1/2}}_{= X^T} \quad \underbrace{\Sigma^{1/2} V}_{= X}$$

## Pattern Analysis (1)

24.05.19

Manifold Learning - (continued)

MDS: "PCA on distances/dissimilarities"

$$C = \left( I - \frac{1}{N} \vec{1} \vec{1}^T \right)$$

$$\text{SVD} \left( \frac{1}{2} C^T D^2 C \right) \Rightarrow \vec{u}^T \Sigma^2 = X$$

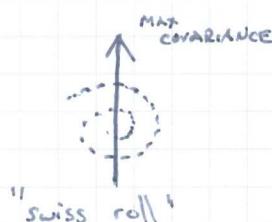
$\uparrow$

zero out  
singular values  
for dim. reduction

Another item to relax from the PCA formulation:

Can we perform a local projection instead of a global one?

Example case (not nicely solvable with PCA/MDS):



[global projection fails]

MDS also has difficulties here:



because the globally applied distances can not distinguish local structures on the manifold.

"ISOMAP": ("Isometric Feature Mapping")

Perform MDS, but on graph distances instead of some global metric like the Euclidean distance.

## Pattern Analysis (2)

24.05.19

Graph Distance:

- 1) Define a local neighborhood, compute Eucl. (or something else) distance to those neighbors  $\Rightarrow$  sparse adjacency matrix
- 2) Compute all-pairs-shortest-paths to obtain all missing distances  
 $\Rightarrow$  This is also called a "geodesic" distance

Locally Linear Embedding (LLE)

Idea: 1) Represent every sample as a linear combination of its neighbors,

2) Search for a lower-dimensional representation with the same or a similar linear combination of neighbors

Note: For the linear constraint  $\sum_{j \in N(x_i)} w_{ij} = 1$  we need a linear programming solver  
 $\Rightarrow$  Let's consider the easier problem with a quadratic constraint  $\sum w_{ij}^2 = 1$   
 $\Rightarrow$  Compute the derivative, solve a linear system

( $N(x_i) \hat{=} \text{"Neighborhood of } x_i\text{"}$ )

1) Search for weights  $w_{ij}$ :

$$\min \sum_i \left\| \vec{x}_i - \sum_{j \in N(x_i)} w_{ij} \vec{x}_j \right\|_2^2$$

subject to  $\sum_{j \in N(\vec{x}_i)} w_{ij} = 1$

2) In the lower-dimensional space  $\mathbb{R}^{d'}$ , find:

Solve for  $\vec{x}'$  in

$$\min \sum_i \left\| \vec{x}'_i - \sum_{j \in N(\vec{x}_i)} w_{ij} \vec{x}'_j \right\|_2^2 \quad \text{subject to } \frac{1}{N} \sum_i \vec{x}'_i \vec{x}'_i^\top = I$$

(covariance is identity)

$$\sum_i \vec{x}'_i = \emptyset \text{ (zero-mean)}$$

## Pattern Analysis (3)

24.05.19

Let's work on the solution of the first part of the problem:

The objective function is invariant to translations:

$$\sum_{i=1}^N \left\| \vec{x}_i - \sum_{j \in N(\vec{x}_i)} w_{ij} \vec{x}_j \right\|_2^2 = \sum_{i=1}^N \left\| \vec{x}_i + \vec{t} - \sum_{j \in N(\vec{x}_i)} w_{ij} (\vec{x}_j + \vec{t}) \right\|_2^2$$

for an arbitrary but fixed  $i$ , we set  $\vec{t} = -\vec{x}_i$

$$\Rightarrow \left\| \vec{x}_i - \vec{x}_i - \sum_j w_{ij} (\vec{x}_j - \vec{x}_i) \right\|_2^2 = \left\| \sum_j w_{ij} (\vec{x}_j - \vec{x}_i) \right\|_2^2$$

$$= \| w_{i1} (\vec{x}_1 - \vec{x}_i) + w_{i2} (\vec{x}_2 - \vec{x}_i) + \dots \|_2^2 = \| M_i w_i \|_2^2 \text{ where}$$

$$\vec{w}_i = \begin{pmatrix} w_{i1} \\ w_{i2} \\ \vdots \end{pmatrix} \text{ and } M_i = \begin{pmatrix} \vec{x}_1 - \vec{x}_i & \vec{x}_2 - \vec{x}_i & \dots \end{pmatrix}$$

$$\Rightarrow \text{minimize } (M_i w_i)^T (M_i w_i) + \underbrace{\lambda(1 - \vec{w}_i^T \vec{w}_i)}_{\text{quadratic constraint}}$$

(Note:  $M_{ij} = 0$  if  $\vec{x}_{ij} \notin N(\vec{x}_i)$ )

$$\frac{\partial}{\partial w_i} (w_i^T M_i^T M_i w_i + \lambda(1 - \vec{w}_i^T \vec{w}_i)) = 2 M_i^T M_i \vec{w}_i - 2 \lambda \vec{w}_i = 0$$

$$\boxed{M_i^T M_i \cdot \vec{w}_i = 2 \lambda \vec{w}_i}$$

Eigenvalue / Eigenvector Problem

( $\Rightarrow$  take eigenvector that belongs to the smallest non-zero eigenvalue)

## Pattern Analysis (4)

24.05.19

2. step: solve second obj function with the calculated weights  $\tilde{w}_i$ ;  
for  $\tilde{x}_i'$  (solution is similar)

### Laplacian Eigenmaps

- 1.) Build an adjacency graph from the feature neighborhoods
- 2.) Compute affinities between neighborhood nodes  $\leftarrow$  weights
- 3.) Perform eigendecomposition of the Graph Laplacian of the weights
- 4.) Low-dim embedding

Step 1:

- Everything within fixed (Euclidean?) distance is a neighbor
- The k-nearest samples are neighbors ( $\Rightarrow$  k-regular graph)

Step 2: Pick any function that decreases with increasing distance

1 specific proposal: heat kernel

$$w_{ij} = e^{-\frac{\|x_i - x_j\|_2^2}{t}}$$

another (non-differentiable) possibility:

binary affinity  $w_{ij} = \begin{cases} 1 & \text{if } x_j \in N(x_i) \\ 0 & \text{else} \end{cases}$

## Pattern Analysis (5)

24.05.19

Step 3, theoretical derivation:

$$\text{minimize} \sum_{i,j} (\vec{x}_i^T - \vec{x}_j^T)^2 \cdot w_{ij}$$

(in the low-dim space)

$$\text{subject to } \vec{x}^T \cdot \vec{x} = 1$$

↑

"Volume constraint"

$$\text{where } D = \begin{pmatrix} \sum_i w_{ii} & & & \\ & \ddots & & \\ & & \sum_j w_{jj} & \\ & & & \ddots \end{pmatrix}$$

Rewrite the objective function:

$$\sum_{i,j} (\vec{x}_i^T - \vec{x}_j^T)^2 \cdot w_{ij} = \sum_{i,j} (\vec{x}_i^T)^2 + \vec{x}_j^T - 2\vec{x}_i^T \vec{x}_j^T \cdot w_{ij}$$

$$= \sum_{i,j} \vec{x}_i^T \cdot w_{ij} + \sum_{i,j} \vec{x}_j^T \cdot w_{ij} - 2 \sum_{i,j} \vec{x}_i^T \vec{x}_j^T w_{ij} =$$

$$= 2 \sum_{i,j} \vec{x}_i^T \cdot w_{ij} - 2 \sum_{i,j} \vec{x}_i^T \vec{x}_j^T w_{ij} = 2 \cdot (\vec{x}^T D \vec{x} - \vec{x}^T W \vec{x})$$

$$= 2 \vec{x}^T (D - W) \vec{x}$$

$\uparrow$

$\sum_i \vec{x}_i^T \cdot \sum_j w_{ij}$

Graph Laplacian  $L$ !

$W = [w_{ij}]$

Pattern Analysis ⑥

24.05.19

$$\text{Minimize } \vec{x}^1 \top L \vec{x}^1 \text{ subject to } \vec{x}^1 \top D \vec{x}^1 = 1$$

$$\Rightarrow \frac{\partial}{\partial \vec{x}} \left( \vec{x}^1 \top L \vec{x}^1 + \lambda (1 - \vec{x}^1 \top D \vec{x}^1) \right)$$

$$= 2L \vec{x}^1 - 2D \vec{x}^1 \stackrel{!}{=} 0$$

$$\Rightarrow D^{-1} L \vec{x}^1 = \lambda \vec{x}^1$$

## Pattern Analysis ①

29.05.19

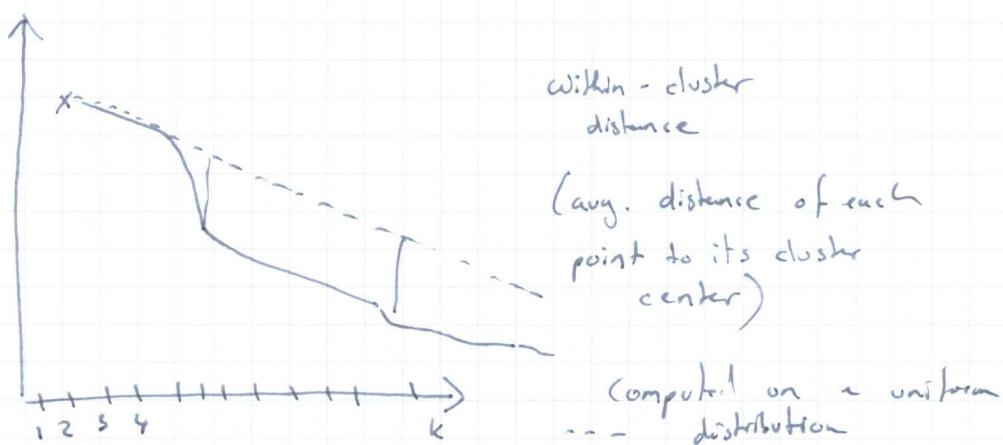
Detour : Gaussian Mixture Models :

$$p(\vec{x}) = \sum_{i=1}^N \beta_i \mathcal{N}(\mu_i, \Sigma_i, \vec{x})$$

↑                      ↑                      ↑  
 Model selection    weighted sum    params that need to be fitted  
 problem:  
 choose  $N$

For clustering - let's say k-means clustering - how did we tackle the model selection problem (i.e. select the number of clusters)?

Gap statistics:



## Pattern Analysis ②

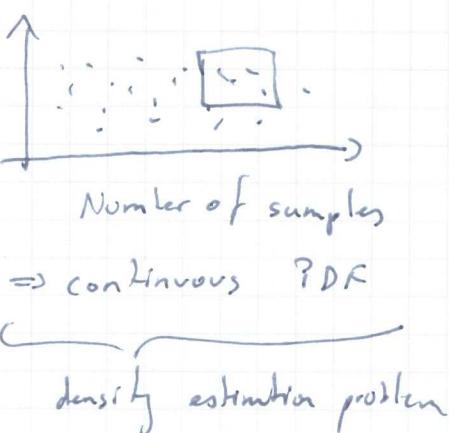
29.05.19

$$\sum k_\lambda (\vec{x}_i - \vec{x})$$

 $\Sigma \dots$ 

"hyper parameter"

← mean shift step

gradient ascent is performed  
with respect to some kernel window

Parzen - Rosenthal - Window

$$p(\vec{x}) = \frac{\sum_{i \in N(\vec{x})} k_\lambda (\vec{x} - \vec{x}_i)}{N \cdot V} \text{ binary}$$

Cross-Validation can be used to  
determine the quality of the chosen  
Parzen window size  $\lambda$  for density  
estimation

## Pattern Analysis (3)

29.05.19

Cross validation:

Remove 1 (or  $k$ ) samples from the input set.For a reasonable set of window sizes  $S = \{l_1, l_2, \dots, l_w\}$ , evaluate
$$1 - \vec{p}_{l_i, x_j}(\vec{x}_j) \leftarrow p(\vec{x})$$
, computed for window  $l_i$  and  
without sample  $\vec{x}_j$ ⇒ repeat for all samples  $\vec{x}_1, \dots, \vec{x}_N$ ⇒ choose the window size  $l_i$  with the smallest interpolation  
error

---

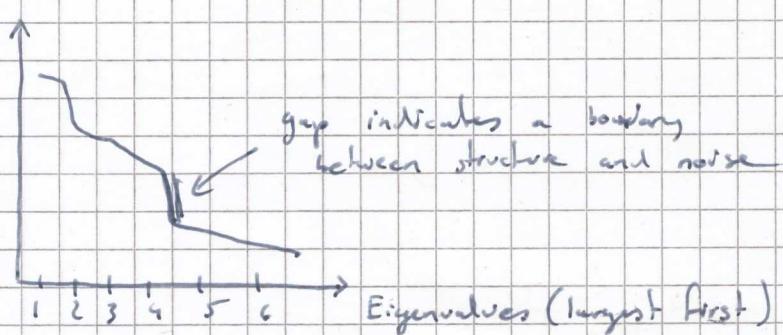
Model selection problem for manifold learning

what is a good (intrinsic!) dimensionality of the data?

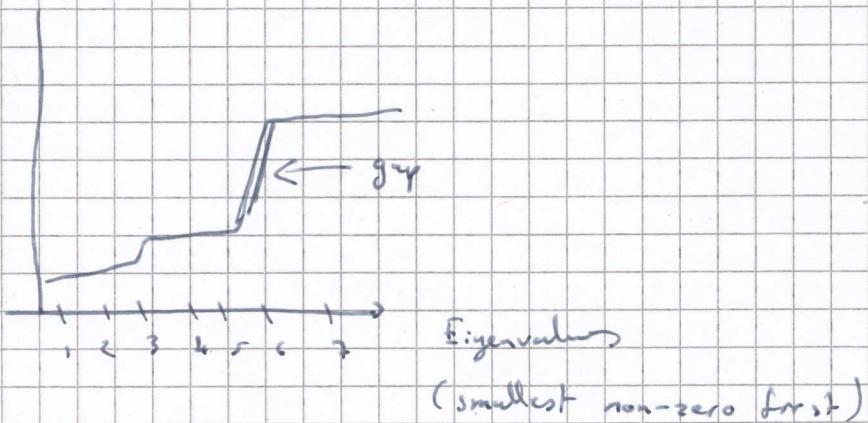
⇒ All methods, that we considered are based on some eigendecomposition.  
The distribution of eigenvalues quantifies the quality of  
the model, similar to the within-cluster distance for  
the k-means clustering.A "good" dimensionality is indicated by a gap in the distribution  
of eigenvalues

For example: pct:

we choose  
the pct  
component that  
contains 99%  
of the  
covariance of  
the data!!



Example LE: (Laplacian Eigenmaps)



### Random Forests

... as a learning-based tool for feature space partitioning



you can use your partitioned feature space for many different tasks, specifically

- classification (1)
- regression (2)
- density estimation (3)
- manifold learning (4)
- semi-supervised learning (5)

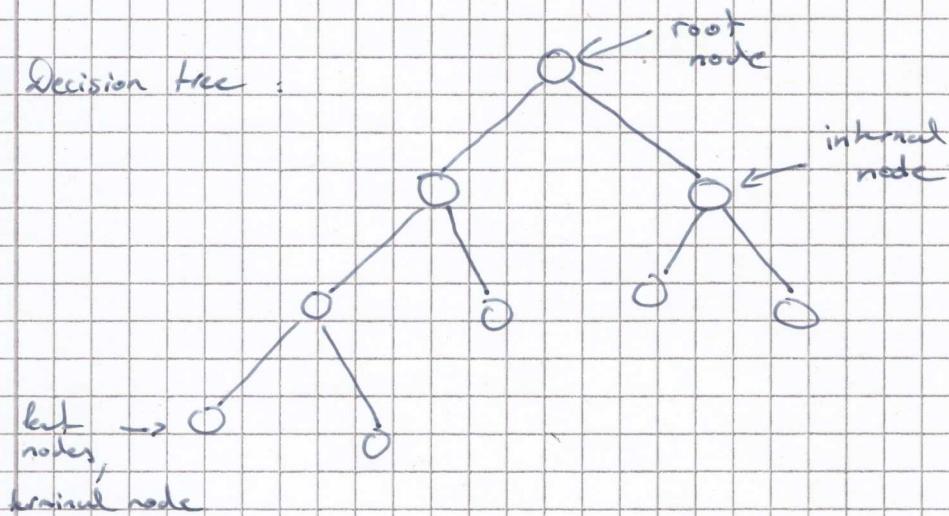
We follow Criminisi, Shotton, Konuklu:

"Decision Forests: A Unified Framework for (1), (2), (3), (4), (5)"

A forest is an ensemble / set of trees.

Decision forest: ensemble of decision trees.

Decision tree:

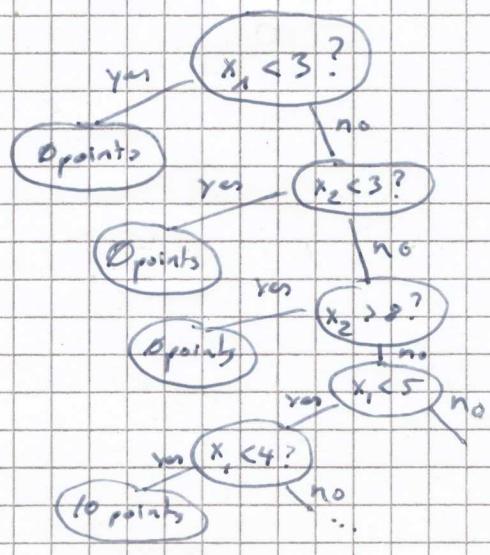
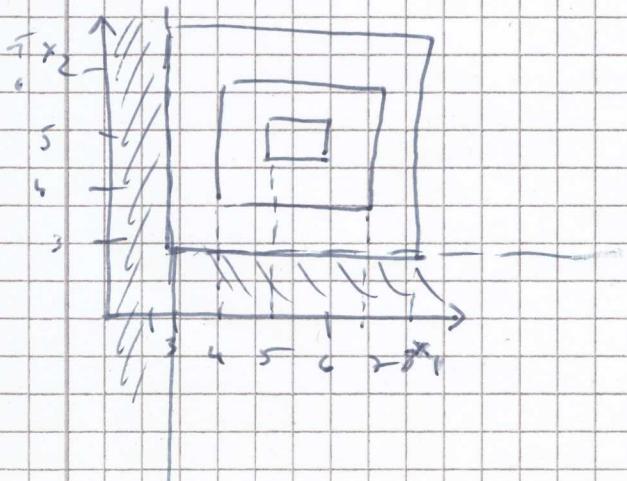


General Workflow:

- start at root node
- answer for each internal node  
with "yes" or "no" to a question
- For "yes" proceed to left child  
For "no" " " "right" "

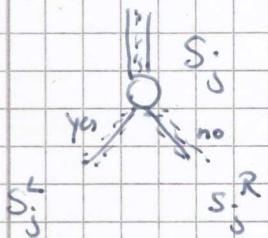
Terminal node gives the result.

"Square cut!"



Random Forests

Training: Determine the functions ("yes/no questions") at each internal node, and the leaf output.



The question for a good node function is: what is a good "splitting function" for all incoming features

⇒ Maximize the Information Gain for our task at hand

$$I(S_j) = H(S_j) - \sum_{i \in \{L, R\}} \frac{|S_j^i|}{|S_j|} \cdot H(S_j^i)$$

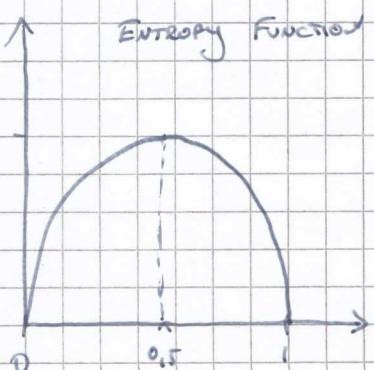
where  $H(S_j)$  denotes the entropy of set  $S_j$

For example, for the case of classification

$$H(S_j) = - \sum_{c \in C} p(c) \cdot \log(p(c))$$

where  $C$  is the set of all class labels, and

$p(c)$  denotes the relative frequency of label  $c$  in set  $S_j$

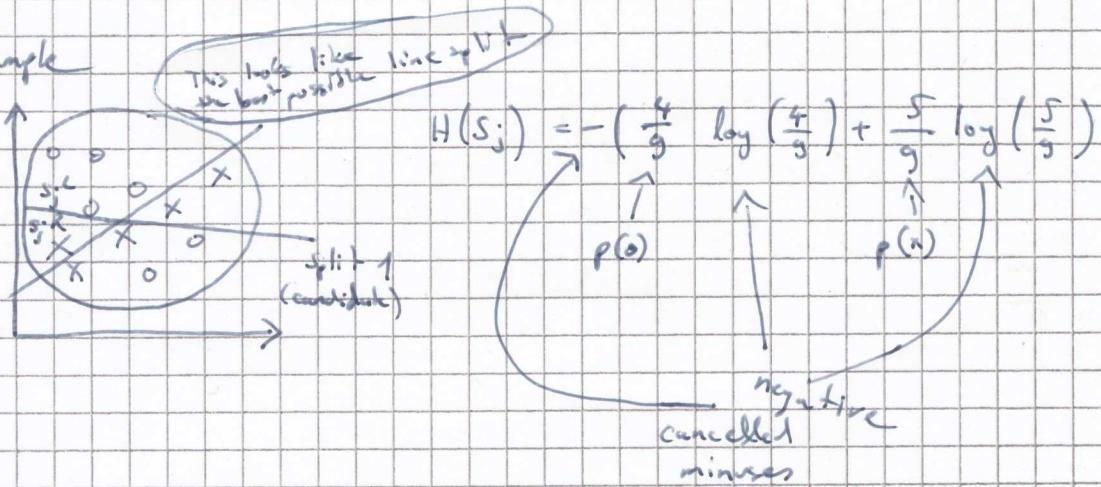


The binary entropy function is lowest if the input is always 0 or always 1  
("not surprising")

Translated to our classification task, the entropy is lowest if all class labels are identical in a subset  $S_j^L$  or  $S_j^R$

$\Rightarrow$  lower entropy in  $S_j^L, S_j^R \Rightarrow$  larger information gain

Example



Let us look for a good split

The entropies of the left and Right subsets of the splits are:

$$H(S_j^L) = - \left( \frac{3}{5} \log \left( \frac{3}{5} \right) + \frac{2}{5} \log \left( \frac{2}{5} \right) \right)$$

$$H(S_j^R) = - \left( \frac{1}{4} \log \left( \frac{1}{4} \right) + \frac{3}{4} \log \left( \frac{3}{4} \right) \right)$$

generally speaking (for classification): we are looking for splitting functions that separate our labelled data as clearly as possible.

Plain Decision Trees are quite prone to overfitting

(consider a greedy selection of always the best information gain)

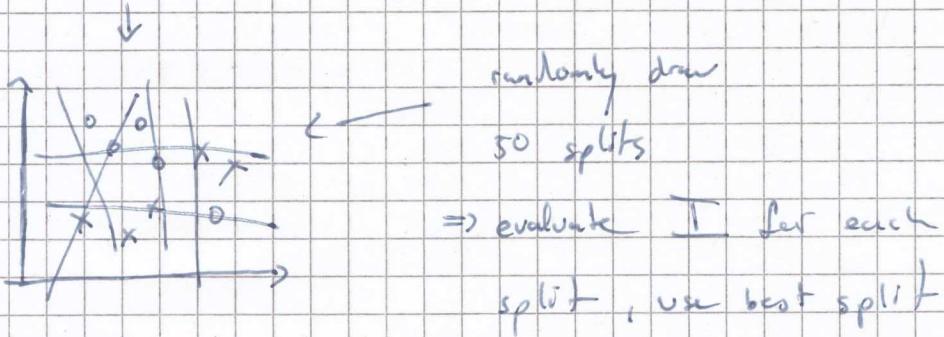
$\Rightarrow$  Breiman's idea: add randomization into the training process

This leads to a "weak learner"  $\Rightarrow$  a tree that is not optimal, but "quite ok" / better than guessing

⇒ Repeat this ~ couple of 100s / 1000s of times to obtain  
a forest of weak decision trees. Average their weak  
opinions ⇒ if the weak learners are statistically independent,  
the ensemble vote converges to the true answer.

Options for randomization:

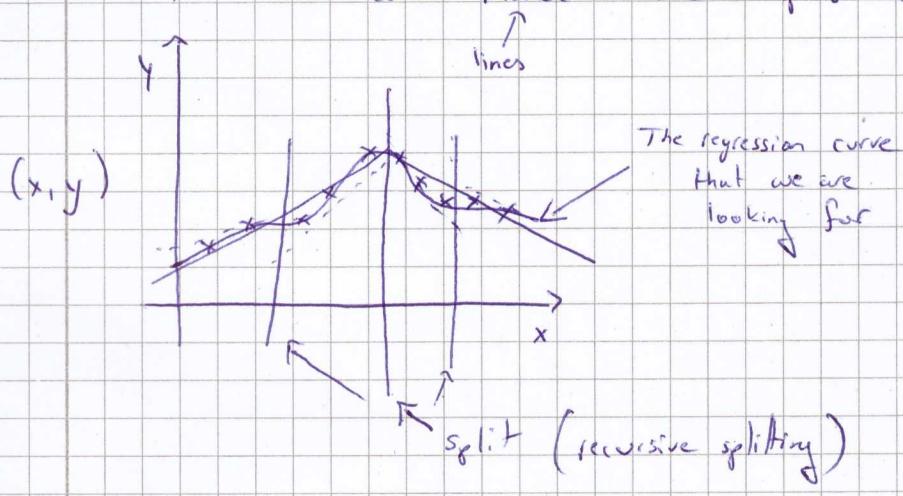
- Randomization is typically set for one decision tree. We can randomize:
  - the subset of samples for training the tree ("bagging")
  - the subset of feature dimensions for the decision functions
  - The parameters of the decision functions



Random Forests

"Regression Forest":

- Redefine information gain for continuous labels  $y$ .
- Fit local curve model to the samples in a leaf node



discontinuities at splitting boundaries

=> Again, by averaging multiple trees, the discontinuities smooth out, and the regression curve is "smooth enough" (for our numerical needs)

Information gain for regression:  $I(s_j) = H(s_j) - \sum_{i \in \{L, R\}} \frac{|s_j^i|}{|s_j|} H(s_j^i)$

$$H(S) = - \frac{1}{|S|} \sum_{x \in S} \int_y p(y|x) \cdot \log(p(y|x))$$

Let  $p(y|x)$  be a Gaussian distribution:

$$p(y|x) = \mathcal{N}(y; \bar{y}(x), \sigma_y^2(x))$$

→ At each position  $x$  we have a Gaussian distribution of possible values (see Bishop's book on how to perform a probabilistic line fit).

## "Density Forest"

⇒ use the Random Forest splitting mechanism to estimate a PDF

- Idea: Fit ~~for each~~<sup>the samples in</sup> leaf node a Gaussian function  
The Gaussian is bounded the boundaries of the leaf (and has discontinuities there)

⇒ Average over many trees to smooth out the discontinuities

Adjust the objective function:  $I(s_j)$  is again identical to classification case

Entropy  $H(s_j)$  for a multivariate Gaussian distribution:

$$H(s) = \frac{1}{2} \cdot \log((2\pi e)^d) \cdot |\Delta(s)|$$

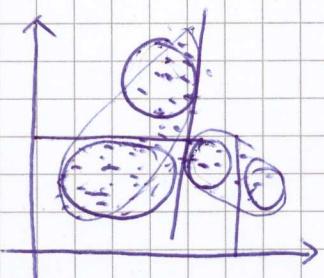
where  $d$  is the dimension of the Gaussian and  $\Delta$  its covariance

→ to evaluate  $H(s)$  you effectively fit a Gaussian to the samples in  $S$ .

Remark: We can plug  $H(s)$  into  $I(s_j)$  and obtain

$$I(s_j) = \log(|\Delta(s_j)|) - \sum_{i \in \{L, R\}} \frac{|s_{j,i}|}{|s_j|} \cdot \log(|\Delta(s_{j,i})|)$$

⇒ Goal: Split to obtain compact Gaussian on each side



"Manifold Forest"

Goal: Manifold Learning

side note:

A parametric model  
has to be finite

Idea: Fit a density forest to the data

use the "uncertainty" of the splitting boundaries

across the trees to determine how tightly two data points are coupled



This is an affinity!

Apply Laplacian Eigenmaps on these affinities

### Affinity Model

Let  $\ell(\vec{v})$  denote the leaf partition function that assigns each feature vector to a leaf.

Let  $k$  denote the number of training samples, then affinity in tree  $t$  between points  $v_i$  and  $v_j$  is

$$w_{ij}^t = e^{-Q^+(v_i, v_j)},$$

$$\text{Let } \vec{d}_{ij} = \vec{v}_i - \vec{v}_j$$

where  $Q^+$  is a distance function.

Choices for  $Q$ :

- Mahalanobis affinity:  $Q^+(v_i, v_j) = \begin{cases} \vec{d}_{ij}^\top (\Lambda_{\ell(v_i)}^+)^{-1} \vec{d}_{ij} & \text{if } \ell(v_i) = \ell(v_j) \\ \infty & \text{otherwise} \end{cases}$

- Binary affinity:

$$Q^+(v_i, v_j) = \begin{cases} 0 & \text{if } \ell(v_i) = \ell(v_j) \\ \infty & \text{otherwise} \end{cases}$$

$$\text{Final affinity } \omega = \frac{1}{T} \cdot \sum_{t=1}^T \omega^t$$

$\Rightarrow$  Compute Laplacian Eigenmaps on  $\omega$

$$L = I - Y^{-\frac{1}{2}} \omega Y^{\frac{1}{2}}, \text{ where } Y = \text{diag}(\sum w_{1j}, \sum w_{2j}, \dots, \sum w_{Nj})$$

Graph Laplacian

$\Rightarrow$  Eigendecomposition

obtain lower-dimensional embedding:

$$L = \begin{pmatrix} 0 & & & \\ & \ddots & & \\ & & \lambda_1 & \\ & & & \ddots & \\ & & & & \lambda_k \end{pmatrix} \quad \text{eigenvector belonging to smallest non-zero eigenvalue}$$

$$v_k = \begin{pmatrix} e_1 \\ e_2 \\ e_3 \end{pmatrix} \quad \text{coordinate vector}$$

$\downarrow$  embedding for  $v_k$  kth-row

Random Forests : A flexible tool to describe sample neighborhoods

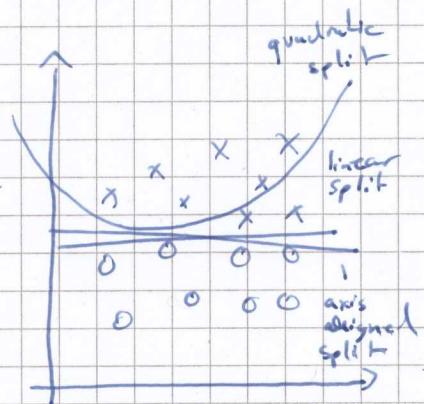
Single Tree: performs a partitioning

Ensemble of trees: "softens" the partitioning from 0-1 labels  
to continuous probabilities.

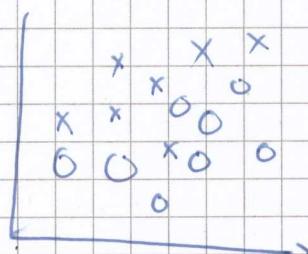
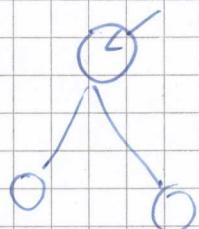


Learned from the data, depending on the tree / forest parameters

- depth of a tree
- number of trees
- early abortion criteria when growing the tree (not discussed)
- Randomization parameters :
  - number of splitting functions to test
  - % of samples to select
  - number of feature dimensions to select} optional
- Parametric form / type of the splitting functions



Tree depth and the number of trees are related when it comes to overfitting / underfitting



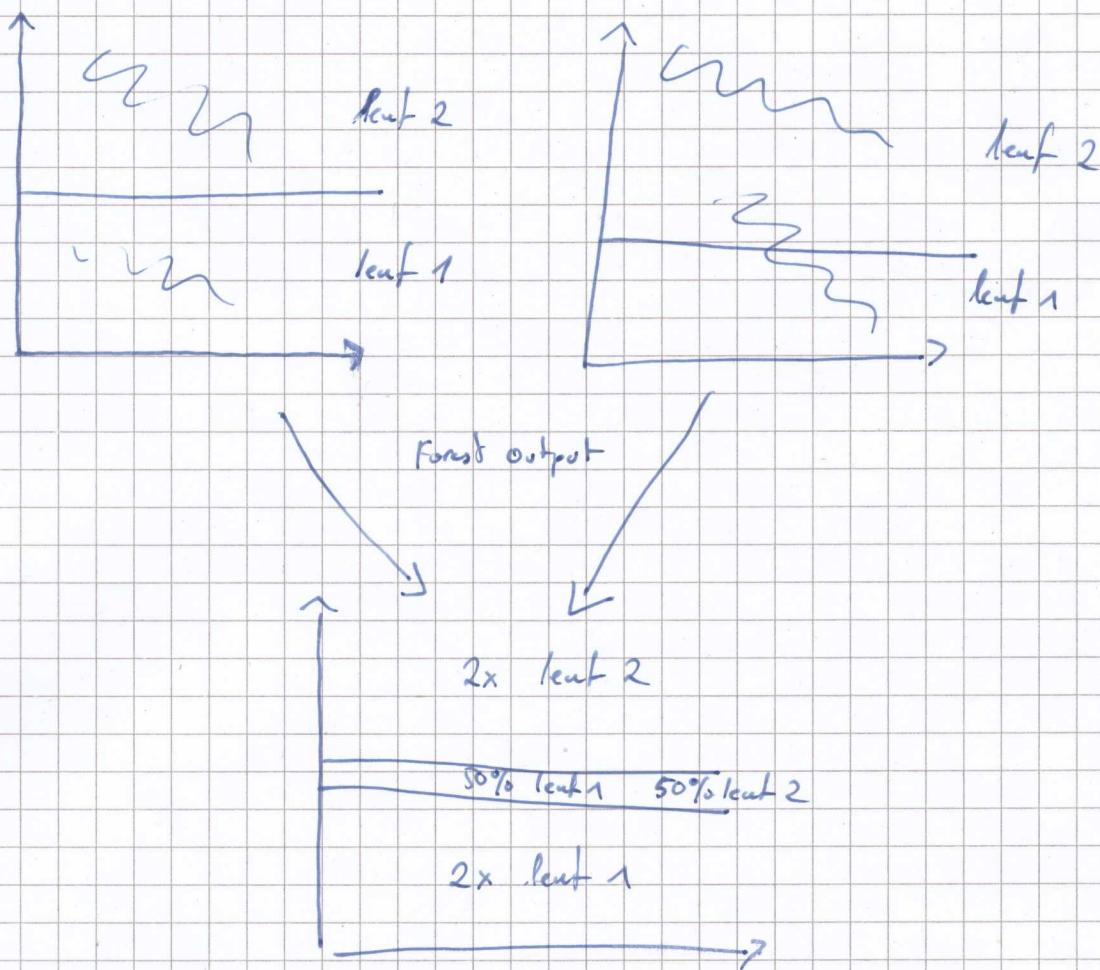
Height 1 : invitation to underfitting (as long as our splitting function is simple - which it always is)

Many trees : blur the decision boundary

→ Height 1 and MAX-NUM of trees maximizes the underfitting

Height  $\log N$  (for  $N$  samples): roughly each sample has its own node → overfit

1 Tree: maximally sharp decision boundary → maximizes overfitting



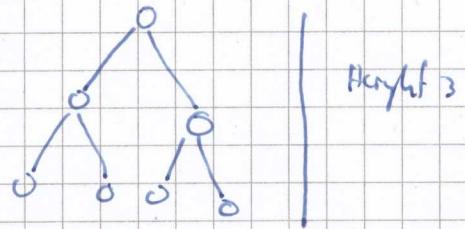
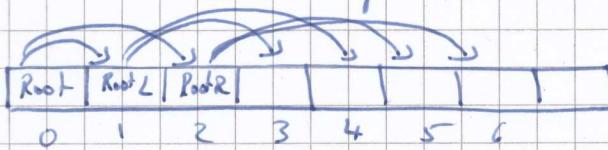
Generally deeper trees also need larger forests to counter overfitting

Good question from audience:

If we have infinite trees, won't the errors introduced by randomization cancel out and lead to overfitting again?

② How to detect overfitting: high training performance, low test performance

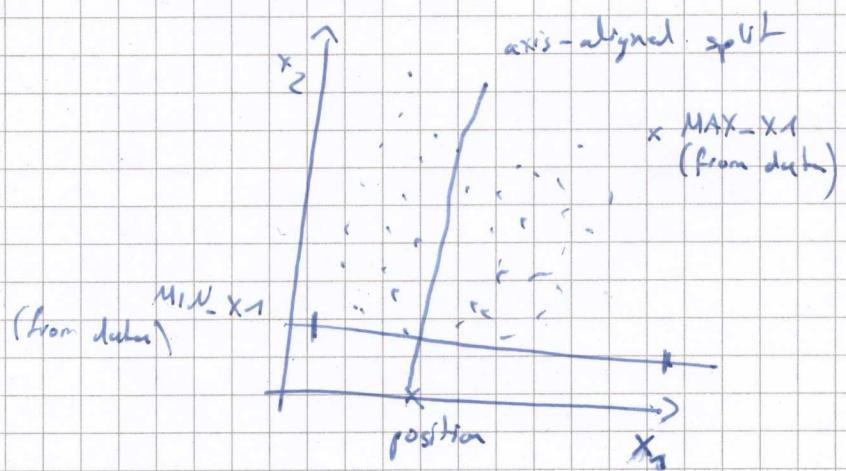
Few notes on the implementation



$2^3 - 1$  nodes

axis-aligned-split

$\left\{ \begin{array}{l} \text{int dimension} \\ \text{double position} \end{array} \right. \quad \left[ \begin{array}{l} x_1/x_c \\ \text{some value} \end{array} \right]$



Sample decision functions:

randomly draw uniformly distributed values

for "dimension" and "position"

within their respective value range

How many decision functions do we need to test in total  
for training a single tree?

Height  $H=3$ , no early stopping

Test per node  $C = 500$  candidate splits

$\Rightarrow$  sample  $2^{(H-1)} \cdot C \cdot 2$  parameters

$$3 \cdot 500 \cdot 2 = 3000$$

$\Rightarrow$  Determine which split is best; (for classification)

for  
each  
at most split

Number of samples  $N$  is constant per tree level

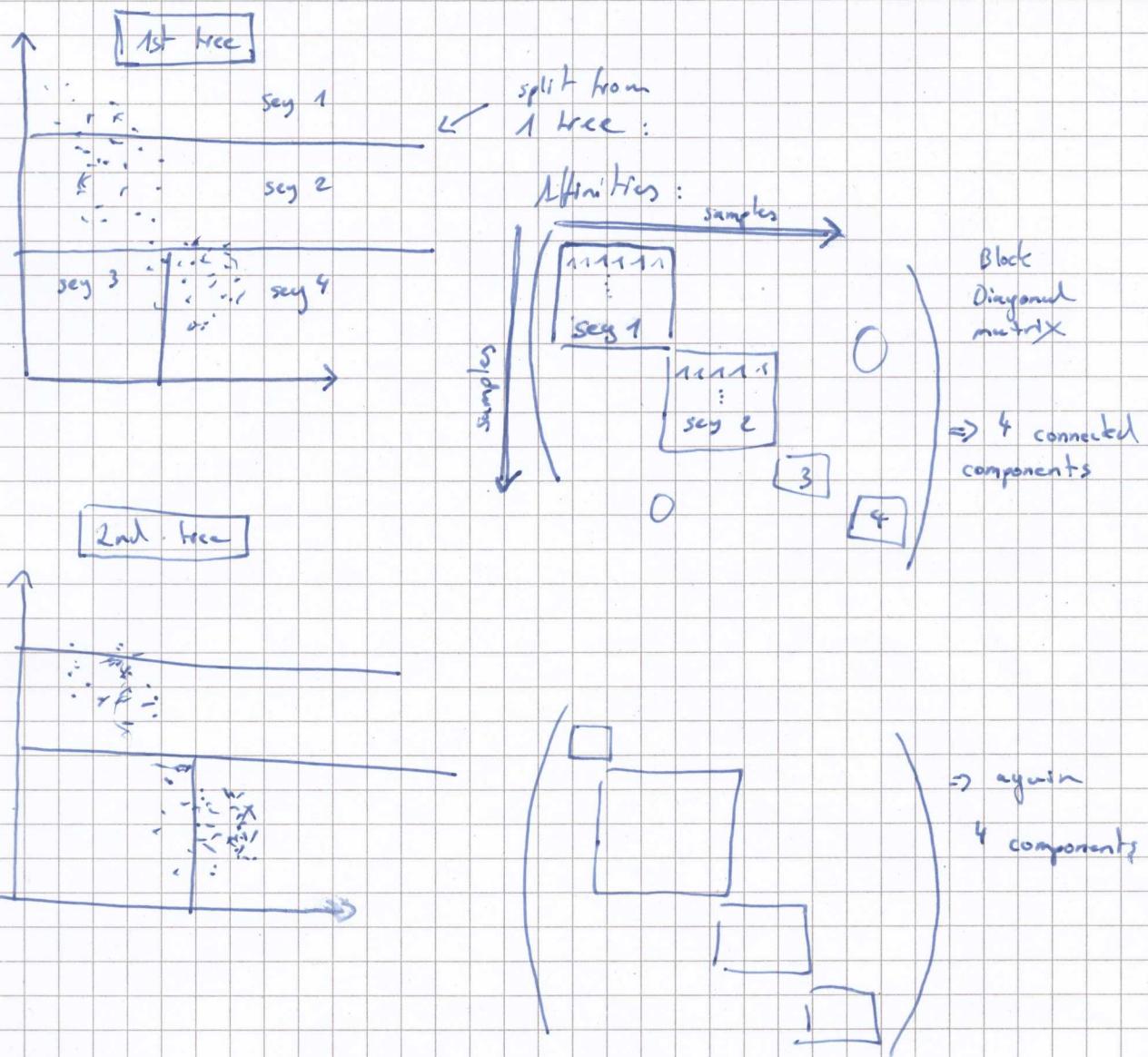
$$N \cdot (H-1) \cdot C = N \cdot 1000$$

$\Rightarrow$  For, say, 5000 trees costs grow by a factor of 5000,  
but each tree is independent  $\Rightarrow$  parallelize : each tree can  
go on each node

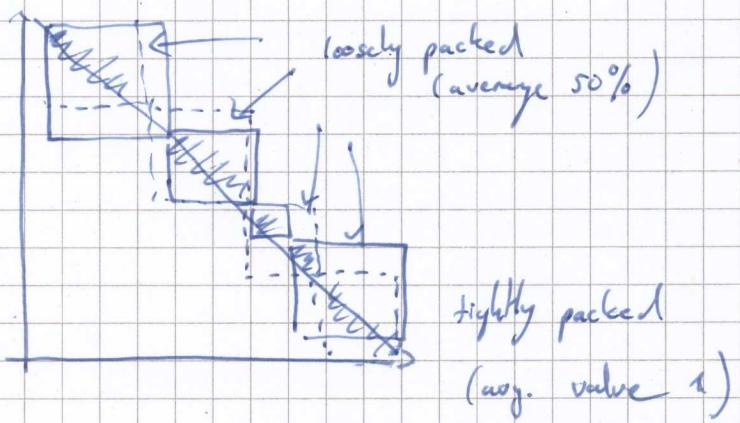
---

A brief sketch on the Manifold Forests :

A manifold forest is Laplacian Eigenmaps that use the partition function of a density forest to define its affinities



Forest output: Average



The loosely packed links can be stretched a little bit when looking for a lower-dimensional manifold.

Nice: 😊 The "rigidness" is deduced from the variations in the tree training

## Probabilistic Graphical Models

- ↳ Hidden Markov Models (HMM)
- ↳ Markov Random Fields (MRF)

Introduction: [ermongroup.github.io / cs228-notes /](https://ermongroup.github.io/cs228-notes/)

Goal: to effectively and efficiently work with a huge joint probability function  $p(x_1, x_2, x_3, \dots, x_n)$

A head-on model for  $p(x_1, \dots, x_n)$  is almost certainly intractable,  
consider for example the space of all spoken sentences  
or the space of all pictures

→ A key component in a probabilistic model are the independence assumptions

$$p(A, B) = p(A) \cdot p(B) \text{ iff } A \text{ and } B \text{ independent}$$

A compact Bayesian network is a distribution in which each factor depends only on a small number of "ancestor variables"  $x_{\pi_i}$ :

$$p(x_1, \dots, x_n) = p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_2, x_1) \cdot \dots \cdot p(x_n | x_{n-1}, x_{n-2}, \dots, x_1)$$

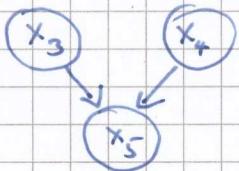
(Chain rule for Bayes formula)

$$p(x_i | x_{i-1}, \dots, x_1) = p(x_i | x_{\pi_i})$$

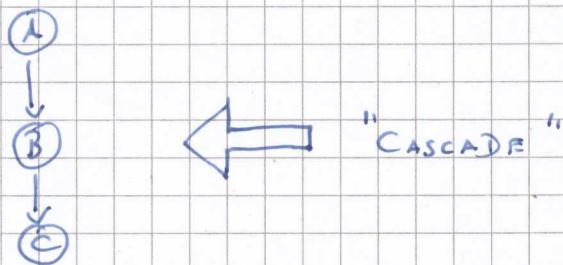
(for example  $x_{A_5} = \{x_3, x_4\}$ )

Graphical representation:

Express a conditional probability as directed edges in a graph, where the variables  $x_1, \dots, x_n$  are the nodes



3-node structures to study independence:

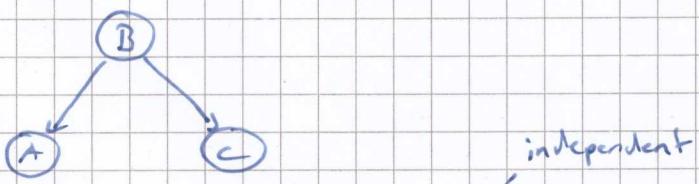


If  $B$  is observed, then  $A \perp\!\!\!\perp C | B$ , that is  $A$  and  $C$  given  $B$  are independent. Conversely, if  $B$  is not observed, then  $A \not\perp\!\!\!\perp C | B$ .

$$p(C|B)$$

$$p(B|A)$$

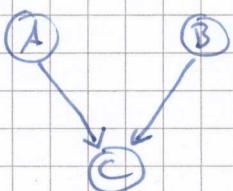
Because  $C$  depends only on  $B$ . If  $B$  is already given, the  $p(B|A)$  doesn't play any role for  $C$  and hence  $A$  and  $C$  are independent

Common Parent

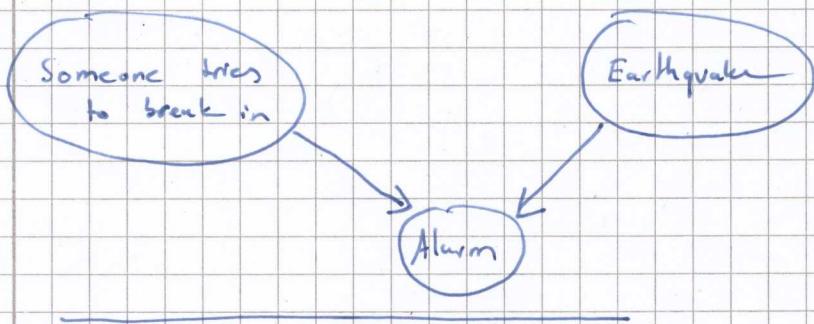
independent

If B is observed, then  $A \perp\!\!\!\perp C \mid B$ . If B is observed  
then  $A \not\perp\!\!\!\perp C$

Because: B holds all information that affects A and C

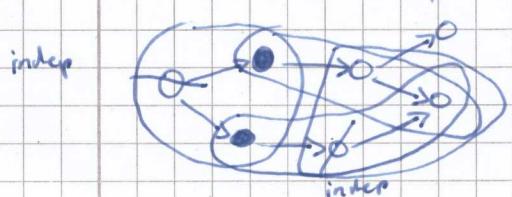
'V'-structure

If we know c, then  $A \not\perp\!\!\!\perp B$ . Conversely, if  
c is unobserved, then  $A \perp\!\!\!\perp B$ . "Explaining away"



Q, W are d-separated when variables O are observed  
if they are not connected by any active path

Active path: A path that consists of triplets, where  
for each triplet the variables are dependent.

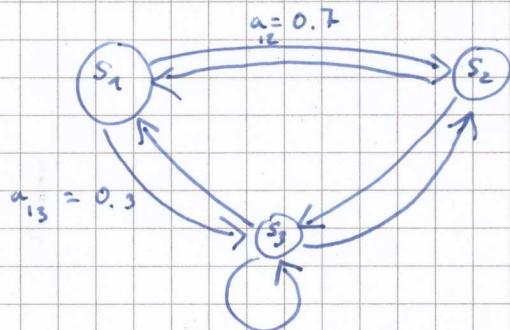


## Hidden Markov Model (HMM)

Probabilistic graphical model. The graph is directed.

$$\begin{aligned} \text{Markov assumption: } & P(q_t = s_i \mid q_{t-1} = s_j, q_{t-2} = s_k, \dots) \\ & = P(q_t = s_i \mid q_{t-1} = s_j) \end{aligned}$$

each conditional probability consists of only one ancestor



nodes in HMM language: "states"

edges: "state transitions"

an edge induces a conditional probability, for example consider

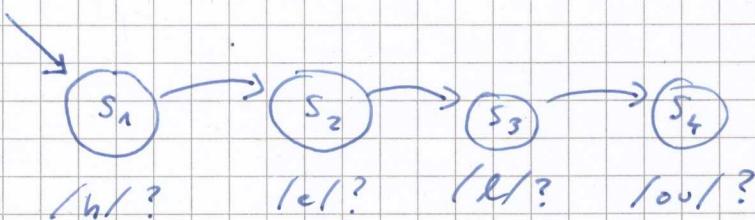
$$P(q_t = s_2 \mid q_{t-1} = s_1)$$

Most popular application: speech processing

More general: sequential data.

Example: Model a word with a HMM

/h/ /e/ /l/ /oo/ : phonemes ("sound atoms")

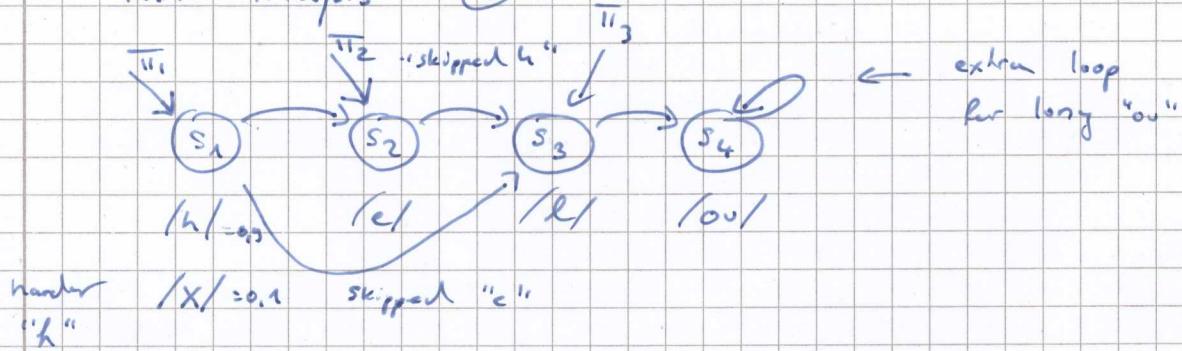


Ambiguities in pronunciation (dialects, accents, slang)

should also be modelled: we introduce additional possibilities for the word of interest (here: "hello") in the graph

## Pattern analysis (2)

19.06.19



← extra loop  
for long "ou"

State transition matrix:  $A$

Output probability matrix:  $B$      $b_i(o_j)$  is the probability  
to observe / produce in  
state  $i$  the symbol  $o_j$ .

$$b_1(o_j = /h/) = 0,9 \quad \sum = 1$$

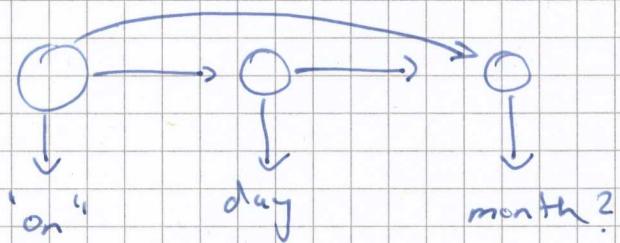
$$b_1(o_j = /x/) = 0,1$$

Starting probabilities  $\pi_1, \dots, \pi_Q$   
max number  
of states

Overall application: Ticket booking system

Construct a dialogue  $\Rightarrow$  prior for the content of the  
next sentence

Absolute date:



"Hidden" is the actual state : We only observe the phonems (or more generally : the output symbols), and are interested in the question how well a HMM matches that sequence, but we don't care about the exact visited states.

$$\begin{aligned} p(\langle o_1, \dots, o_T \rangle, \langle s_1, \dots, s_Q \rangle) &= \\ p(\langle s_1, \dots, s_Q \rangle) \cdot p(\langle o_1, \dots, o_T \rangle | \langle s_1, \dots, s_Q \rangle) &= \underline{\underline{p(s_i | s_{i-1})}} \end{aligned}$$

Markov property: Joint probabilities can be factorized into small terms, because every probability conditioned on state depends only on the last state



Task 1: How to recognize a word?

Let  $\lambda = (A, B, \pi)$  denote our HMM for a specific word (e.g. "hello")

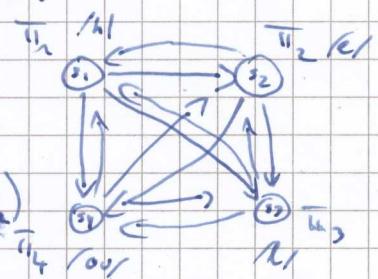
$\Rightarrow$  we seek  $p(\langle o_1, \dots, o_T \rangle | \lambda_{\text{hello}})$

$\Rightarrow$  We need to marginalize out all hidden state sequences  $\langle s_1, \dots, s_Q \rangle$

Brute force: Start at  $\pi_i^1$

$$\sum_{i_1=1}^{Q=4} \pi_{i_1}^1 \cdot b_{i_1}(o_1) \cdot \sum_{i_2=1}^Q a_{i_1 i_2} \cdot b_{i_2}(o_2) \cdots$$

$i_1$



$$b_i(o) = \begin{pmatrix} 0.61 \\ 0.09 \\ 0.08 \\ 0.9 \end{pmatrix} \begin{matrix} /h/ \\ /t/ \\ /l/ \\ /ou/ \end{matrix}$$

$$\cdot \left( \sum_{i_2=1}^Q a_{i_2 i_3} \cdot b_{i_3}(o_{j_3}) \right) \cdot \left( \sum_{i_4=1}^Q a_{i_3 i_4} \cdot s_{i_4}(o_{j_4}) \right)$$

## Dynamic Programming

Due to Markov property  
we only need to cache  
the probabilities of  
the previous state

		State $i_3$
		1 ... Q
State $i_4$		1
1	$a_{11}$	
1	$a_{12}$	
1	$a_{13}$	0 0 0
1	$a_{14}$	$a_{21} a_{22} a_{34} a_{41}$

0.65

Create a list of  
state numbers and  
probabilities for producing  
 $o_{j_4}$  (1001)

- $s_1$  [0.0]
- $s_2$  [0.5]
- $s_3$  [0.6]
- $s_4$  [0.9]

HMM

Given: sequence of samples / observations

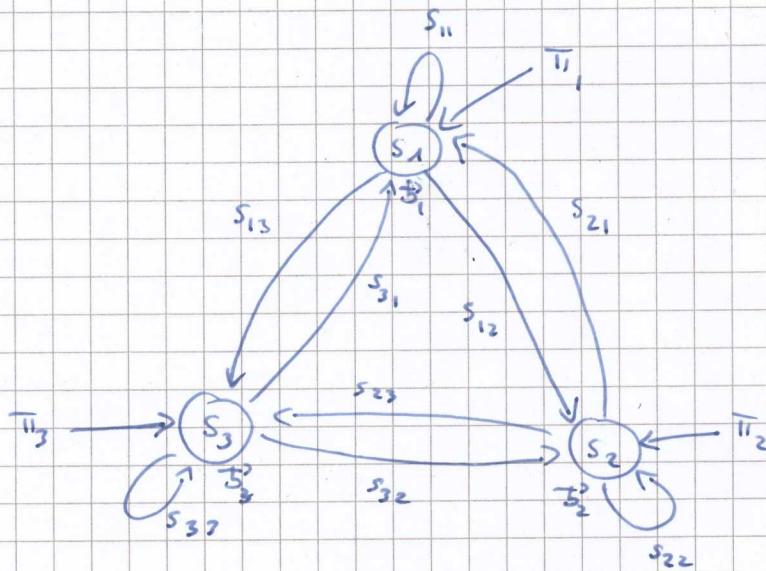
$$\langle o_1, \dots, o_n \rangle$$

Desired goal: model for representing  $\langle o_1, \dots, o_n \rangle$ , for example to compute a probability, that  $\langle o_1, \dots, o_n \rangle$  describes a certain label (speech: a word from the vocabulary)

Rabiner: "A Tutorial on Hidden Markov Models with Selected..."

Approach: Use a probabilistic model, expressed in a directed graph.

3-state example sketch:



Length of input sequence: N

# of states: M

# of symbols in alphabet: V

$\pi$ : PDF of starting probabilities

$A = [a_{s_i s_j}]$  Transition probabilities

set of states  $S = \{s_1, \dots, s_M\}$

$$A \in \mathbb{R}^{M \times M}$$

$B = [b_{s_i}(o_j)]$  Output probabilities

$$B \in \mathbb{R}^{M \times V}$$

"Hidden" is the sequence of states for matching of  $\langle o_1, \dots, o_n \rangle$

The HMM models the joint probability  $p(\langle o_1, \dots, o_n \rangle, \langle s_1, \dots, s_n \rangle)$

The Markov property implies that  $p(s; \text{at time } t | s_j \text{ at time } t-1)$

$$= p(s; \text{at time } t | s_j(t-1), s_{j-1}(t-2), \dots, s_1(1))$$

↑  
"at time" t-1 >

that is, the current state only depends on the previous state  
and not on the whole previous path (or future path)

$\Rightarrow$  this makes inference tractable

### 3 Methods to work with HMMs

1.) Matching: find  $p(o_1, \dots, o_n | \lambda)$ : Forward Algorithm

Backward Algorithm

$$\lambda = (A, B, \pi)$$

2.) Most likely state sequence for a given observation

$$\underset{\langle s_1, \dots, s_n \rangle}{\operatorname{argmax}} \quad p(o_1, \dots, o_n, \langle s_1, \dots, s_n \rangle)$$

Viterbi algorithm

3.) Training: find  $\lambda = (A, B, \pi)$  given some training samples

Baum - Welch - Formula(e)

Task 1.)  $p(o_1, \dots, o_n, \langle s_1, \dots, s_n \rangle) =$

$$\underbrace{p(s_1, \dots, s_n)}_{\substack{\text{constant for} \\ \text{different} \\ \text{inputs}}} \cdot \underbrace{p(o_1, \dots, o_n | s_1, \dots, s_n)}_{\text{goal: } p(o_1, \dots, o_n)}$$

$$p(o_1, \dots, o_n) = \sum_{s_1=1}^M \sum_{s_2=1}^M \sum_{s_3=1}^M \dots \sum_{s_n=1}^M \left( \pi_{s_1} \cdot b_{s_1}(o_1) \cdot a_{s_1 s_2} \cdot b_{s_2}(o_2) \cdot \dots \right.$$

11111111111111  
N-digit number  
with base M

$$\left. \dots a_{s_{n-1} s_n} \cdot b_{s_n}(o_n) \right)$$

Naive Approach is too expensive (exponential), but there is a dynamic programming solution

$$\sum_{s_1=1}^M \pi_{s_1} b_{s_1}(o_1) \cdot \left( \sum_{s_2=1}^M a_{s_1 s_2} \cdot b_{s_2}(o_2) \cdot \left( \sum_{s_3=1}^M \dots \left( \sum_{s_n=1}^M a_{s_{n-1} s_n} \cdot b_{s_n}(o_n) \right) \right) \right)$$

### FORWARD ALGORITHM

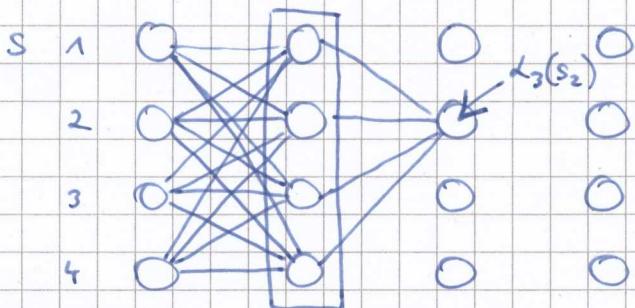
$$1.) \quad t=1 : \quad \alpha_1(s_1) = \pi_{s_1} \cdot b_{s_1}(o_1)$$

$$2.) \quad 1 < t \leq N : \quad \alpha_t(s_t) = \sum_{s_{t-1}=1}^M \alpha_{t-1}(s_{t-1}) \cdot a_{s_{t-1} s_t} \cdot b_{s_t}(o_t)$$

$$3.) \text{ afterwards} \quad p(o_1, \dots, o_N) = \sum_{s_N=1}^M \alpha_N(s_N)$$

$\alpha_t(s_t)$  caches the probability for all paths starting at time 1 up until time  $t$ , that meet at time  $t$  in state  $s_t$ :

$$t = 1 \quad 2 \quad 3 \quad 4 \quad \dots \quad N$$



### BACKWARD ALGORITHM

$$1.) \quad \beta_N = 1$$

$$2.) \quad 1 \leq t \leq N : \quad \beta_t(s_t) = \sum_{s_{t+1}=1}^M a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{t+1}(s_{t+1})$$

$$3.) \quad t=1 : \quad \sum_{s_1=1}^M \pi_{s_1} \cdot b_{s_1}(o_1) \cdot \beta_1(s_1) \quad \begin{array}{l} (\text{wikipedia?}) \\ \text{Termination?} \end{array}$$

It's missing in the paper ...

$$p(o_1, \dots, o_n) = \sum_{i=1}^M \overline{\pi_i} \beta_i(s_i) \cdot b_i(o_n)$$

(↑ wrong?)

For exam: read through the "must-read" papers

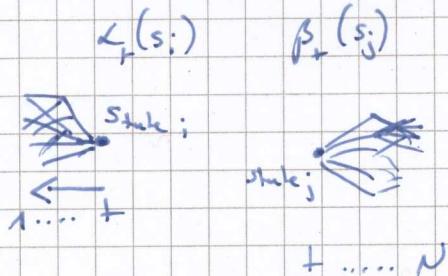
- write down equations and explain

### HMM (continued)

Three methods:

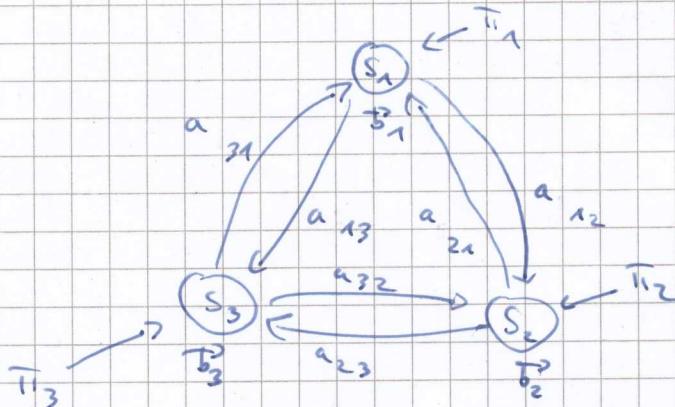
$$\begin{aligned} 1.) \quad & p(o_1, \dots, o_n | \lambda) \\ & \uparrow \\ & \text{or} \\ & p(o_1, \dots, o_n) \end{aligned}$$

: FORWARD - / BACKWARD algorithm



2.) Most likely state sequence for  $\langle o_1, \dots, o_n \rangle$  : Viterbi

3.) Training of  $\lambda = (\Lambda, B, \pi)$ : Baum-Welch - formulas .



Example path

$$\pi_2 \cdot b_2(o_n) \cdot a_{21} \cdot b_1(o_1)$$

Missing Termination step in BACKWARD: (Footnote in paper)

omitted b.c. we never backtrack

to the beginning (since we use FORWARD)  
from there, I guess)

Solution to Task 2: (Viterbi algorithm):

Most likely state sequence:

Method is identical to FORWARD algorithm, with two exceptions:

1.) Replace every sum by max



2.) Technical detail: store the actual path in a second list  $\gamma$ .

Solution to Task 3 (Training):

Determine  $A, B, \pi$ : no analytical solution  $\Rightarrow$  solve iteratively

Auxiliary quantities

$$\alpha_t(s_i, s_j) = p(q_t = s_i, q_{t+1} = s_j)$$

↑  
probability for  
transition  
 $s_i \rightarrow s_j$  at  $t$

$$= \alpha_t(s_i) \cdot a_{s_i s_j} \cdot b_{s_j}(o_{t+1}) \cdot \beta_{t+1}(s_j)$$

$$\underbrace{\sum_{k=1}^M \sum_{l=1}^M \alpha_t(s_k) \cdot a_{s_k s_l} \cdot b_{s_l}(o_{t+1}) \cdot \beta_{t+1}(s_l)}$$

Normalization for all transitions

between time  $t$  and  $t+1$

$A = [a_{ij}]$  matrix of all transition probabilities

$B = [b_i(o_j)]$  matrix of all output probabilities

$\pi$ : starting probabilities

(given: - example sequences

$\langle o_1, \dots, o_n \rangle$  ("training data")

- number of states

- rough connection structure

(e.g. complete graph,  
left-right HMM)

↑  
triangle transition  
matrix

$$\gamma_+(s_i) = \sum_{s_j=1}^M \xi_+(s_i, s_j)$$

↑  
probability of being  
in state  $s_i$  at  
time  $t$

---

$$\sum_{t=1}^T \xi_+(s_i, s_j) : \text{expected \# of transitions } s_i \rightarrow s_j$$

$$\sum_{t=1}^T \xi_-(s_i) : \text{expected \# of transitions from } s_i$$

### Baum-Welch Formulas:

EM-style-iteration:

↑  
Expectation Maximization

Iterate {

$\bar{\pi}_i = \text{expected \# of times in } s_i \text{ at time } t=1$

$$= \gamma_+(s_i)$$

$$\bar{a}_{ij} = \frac{\text{expected \# of transitions } s_i \rightarrow s_j}{\text{expected \# of transitions from } s_i}$$

$$= \frac{\sum_{t=1}^T \xi_+(s_i, s_j)}{\sum_{t=1}^T \xi_-(s_i)}$$

$$t_{s_j}(k) = \frac{\text{expected \# of times in } s_j \text{ and observing symbol } v_k}{\text{expected \# of times in } s_j}$$

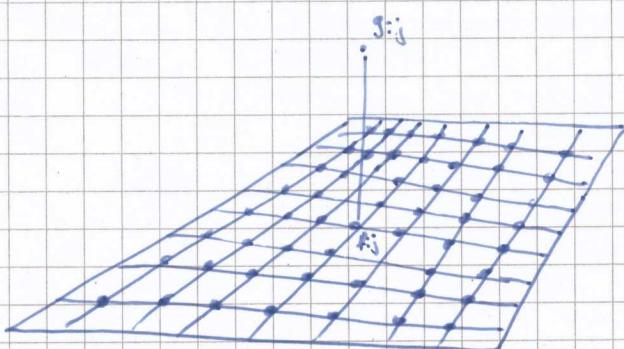
$$= \frac{\sum_{t=1}^T \xi_+(s_j, o_t) \text{ where } o_t = v_k}{\sum_{t=1}^T \xi_+(s_j)}$$

Markov Random Field

Probabilistic graphical model on undirected graph

The graph structure is not limited to modeling 1-D sequences (like in an HMM), but can cover arbitrary neighborhood relations.

Important special case: consider the pixel grid of an image as a graph:



Each hidden variable depends on:

- the associated observation
- its neighborhood of other hidden variables

It is independent of all other quantities  $\Rightarrow$  "the trick" to make inference over

$$p(f_{ij}, g_{ij}) \text{ tractable}$$

↑  
all hidden vars

↑  
all observations

Observations: grid of noisy pixels  
 $[g_{ij}]$

Hidden Variables: grid of ideal (noiseless) image pixels  
 $[f_{ij}]$

More specifically, we model the joint probability  $p(f_{11}, f_{12}, \dots, f_{m1}, g_{11}, \dots, g_{m1})$

$$p([f_{ij}], [g_{ij}]) = p([g_{ij}] | [f_{ij}]) \cdot p([f_{ij}])$$

The Markov assumption limits the statistical dependency of the hidden variables:  $p([f_{ij}]) =$

$$= p(f_{ij} | f_{11}, f_{12}, \dots, f_{ij-1}, f_{ij+1}, \dots, f_{m1}) \cdot \\ \circ p(f_{11}, f_{12}, \dots, f_{ij-1}, f_{ij+1}, \dots, f_{m1})$$

Markov:

$$p(f_{ij} | f_{11}, f_{12}, \dots, f_{ij-1}, f_{ij+1}, \dots, f_{m1}) = p(f_{ij} | N(f_{ij}))$$

↑  
neighborhood of  $f_{ij}$

$$\text{for example: } N(f_{ij}) = \{f_{ij-1}, f_{i-1,j}, f_{ij+1}, f_{i+1,j}\}$$

→ goal is to look for the most probable assignment of values to the hidden variables

$$[f_{ij}] = \underset{f_{ij}}{\operatorname{argmax}} \ p([g_{ij}] | [f_{ij}]) \cdot p([f_{ij}])$$

2 open problems: 1.) How to define  $p([g_{ij}] | [f_{ij}])$  and  $p([f_{ij}])$ ?

2.) How can we efficiently find the  $\operatorname{argmax}$ ?

Task 1:

$p([g_{ij}] | [f_{ij}])$  defines the relationship of a hidden variable to one or more observations.

For the example of denoising, we could define this as a PDF for, e.g. Gaussian-distributed noise in the observations  $[g_{ij}]$ .

This term is also called "data term"



describes relationship to the data / observations

$p([f_{ij}])$  defines the relationship of hidden variables within a local neighborhood.

For the denoising example, this PDF can capture the fact, that the ideal image  $[f_{ij}]$  is smooth, i.e., neighboring pixels / labels are very often identical.

conceptually, we are done. However for most tasks it is very difficult to find "good" pdfs that translate our concept into the language of probabilistic modeling.

⇒ Solution: The Hammersley-Clifford Theorem states, that a MRF is (essentially) equivalent to a Gibbs Random Field (GRF).

This allows to cast our concept into Gibbs Potentials, which are much more intuitive

A GRF is given by

$$p(\vec{x}) = \frac{1}{Z} \cdot e^{-H(\vec{x})}$$

where  $Z = \sum_{\vec{x}'} H(\vec{x}')$  is called "Partition Function" and

normalizes the probability. It is hard to compute, but doesn't depend on  $\vec{x}$ , so for the maximization of  $\vec{x}$  it is ignored.

$H(\vec{x})$  is an "energy function", consisting of potentials  $V_m(\vec{x})$ :

$$H(\vec{x}) = \sum_{m \in S} V_m(\vec{x})$$

If we minimize the energy function, we maximize the GRF

$\Rightarrow$  Let's define the potentials in such a way, that they are closer to  $\emptyset$  in cases where  $p([\vec{f}_{ij}], [\vec{g}_{ij}])$  is closer to our expectations

Specific example for image smoothing:

$$p([\vec{g}_{ij}] | [\vec{f}_{ij}]) = \prod_{ij} \frac{1}{\sqrt{2\pi \sigma_{ij}^2}} \cdot e^{-\frac{1}{2\sigma_{ij}^2} \cdot (\vec{f}_{ij} - \vec{g}_{ij})^2}$$

$$= \frac{1}{\prod_{ij} \sqrt{2\pi \sigma_{ij}^2}} \cdot e^{-\sum_{ij} \frac{1}{2\sigma_{ij}^2} \cdot (\vec{f}_{ij} - \vec{g}_{ij})^2} \quad \text{a union centered around } \vec{f}_{ij}$$

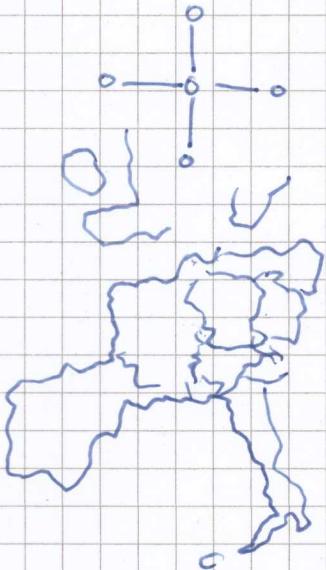
$$p([\vec{f}_{ij}]) = \frac{1}{Z} \cdot e^{-\sum_{ij} \|\nabla f_{ij}\|^2}$$

↑  
potential

## MRF: Markov Random Field

Probabilistic graphical model.

Markov property: probability of a node labeling  
only depends on the neighbors of the node.



Probabilistic Model: MRF

↑  
Hammersley-Clifford-Theorem

Gibbs Random Field (GRF)

GRF: potential functions

Image smoothing example: "data term" / unary potential

$$-\log \left( \text{crf}(g_{ij}, f_{ij}, \theta) \right) = \frac{(f_{ij} - g_{ij})^2}{\sigma}$$

"pairwise potential" / "regularizer"

$$E(f_{ij}, f' \in N(f_{ij})) = \|f_{ij} - f'\|_2^2$$

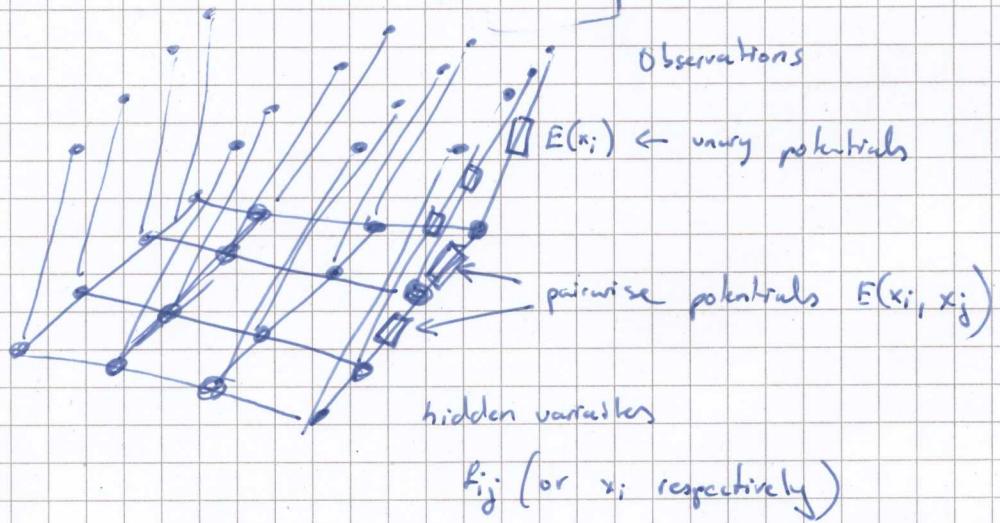
These potentials form an energy function. Minimizing these Gibbs potentials leads (under some mild assumptions) to the same result as maximizing an MRF

If we only consider the energy terms, the unary potentials are  $E(x_i)$ ,

where  $x_i$  denotes the  $i$ -th hidden variable and the pairwise potentials are  $E(x_i, x_j)$ .

The minimization problem is just  $\sum_{i=1}^N E(x_i) + \sum_{i=1}^N \sum_{j=1}^N E(x_i, x_j)$

$$\rho(\dots) = e^{-\sum E(x_i) + \sum \sum E(x_i, x_j)}$$



An efficient solution for a finite set of labels can be computed via graph cuts

Kolmogorov, Zabih: "What Energy Functions Can Be Minimized via Graph Cuts?", PAMI 2004

$\alpha$ -Expansion algorithm for L labels:

while (changes) :

for each label l:

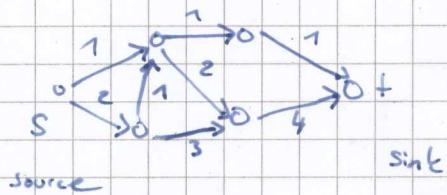
binaryGraphCut(l); // 1vs.all graph cut:

does the number of assignments to label l grow?  
⇒ keep solution

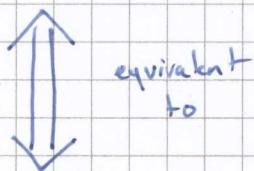
$\Rightarrow$  the core of the label assignment is a binary labeling via graph cuts. This is what we discuss below.

### Maximum flow

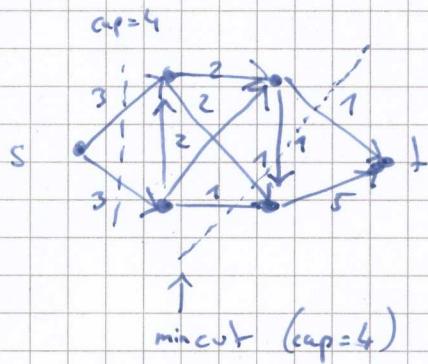
Edge weights = pipe capacities



Question: what is the maximum throughput?



what is the minimum cut?



Separates source  $s$  and sink  $t$ , such that the sum of edges from source to sink is minimal.

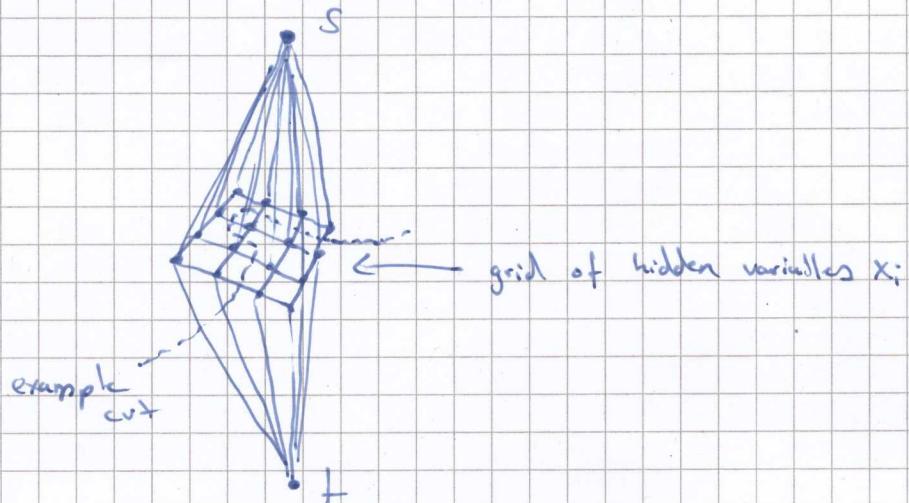
Here:  $\text{MinCut} = 4$

Idea: Identify the nodes in a graph with the hidden variables

$x_i$ : identify source  $s$  with label  $\emptyset$ , drain  $t$  with label 1

Task: Define the graph topology and weights such that the mincut separates the nodes into  $\emptyset$ 's and 1's with an optimal solution to the energy minimization problem

Graph layout:

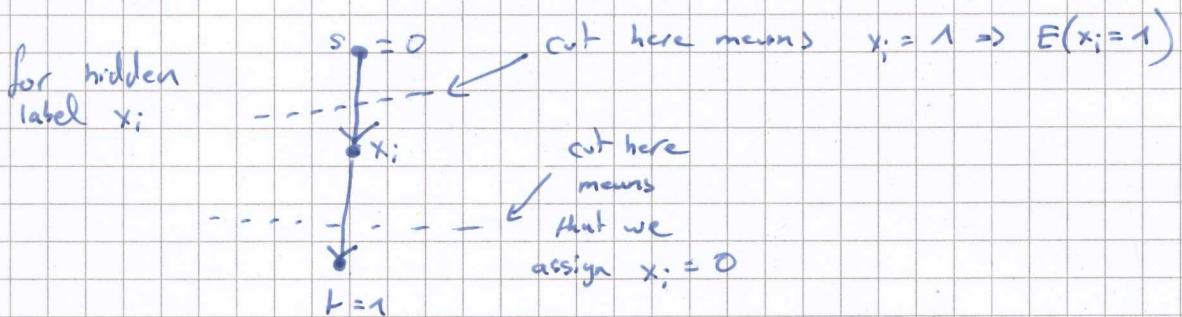


Assignment of edge weights

1.) MinCut works on parts of a graph the same way as the sum of the parts

$\Rightarrow$  decompose  $\sum E(x_i) + \sum \sum E(x_i, x_j)$  into individual terms  
 $E(x_1), E(x_2), \dots, E(x_1, x_2), \dots, E(x_1, x_3)$

2.) Part of the graph for  $E(x_i) \leftarrow$  unary potential



$\Rightarrow$  the cost should be

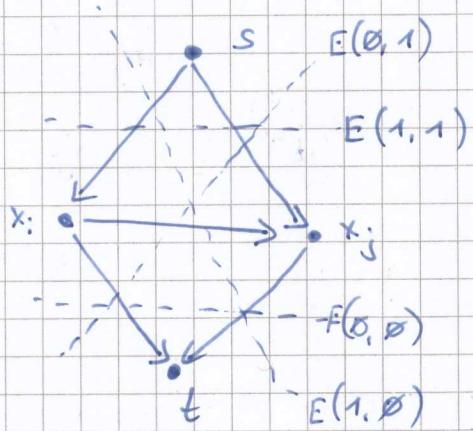
$$E(x_i = 0)$$

(observation)

$$\frac{\rho}{2} \|x_i - g_{ij}\|_2^2$$

For example with our image smoothing function  $E(x_i) = \frac{\rho}{2} \|x_i - g_{ij}\|_2^2$

3.) Part of the graph for  $E(x_i, x_j) \leftarrow$  pairwise potential



Depending on the cut, we can have 4 different label assignments

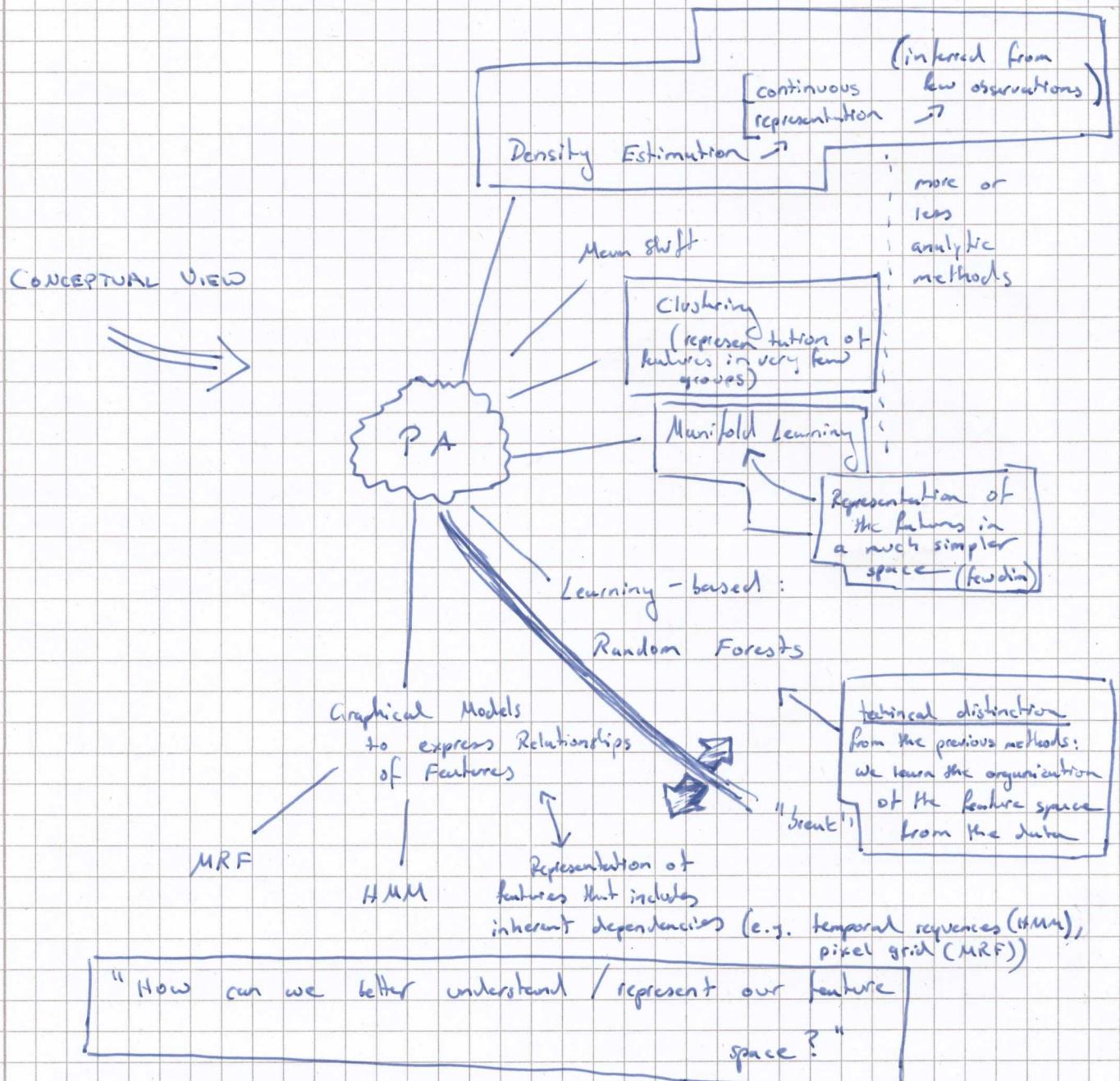
	$x_j$	0	1
$x_i$	0	$E(0,0)$	$E(0,1)$
	1	$E(1,0)$	$E(1,1)$

The conversion CRF  $\rightarrow$  MinCut can always be done (and the MinCut solution is always equivalent to the CRF)  
if the so-called submodularity condition holds:

$$E(\emptyset, \emptyset) + E(1,1) \leq E(0,1) + E(1,0)$$

(assigning identical labels to neighbors may not be more expensive / "worse" than assigning different labels)

$$\left( \|x_i - x_j\| \Rightarrow \text{(holds for smoothing, for example)} \right)$$

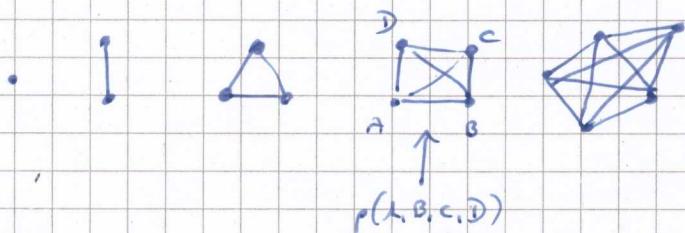
Summary of the class

Technical view: Several tools reappear from time to time.

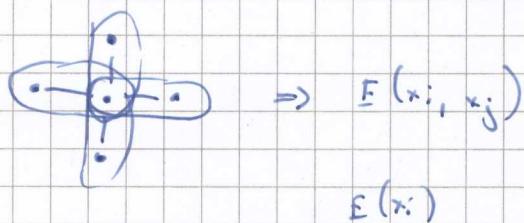
- all methods in the manifold learning part in one way or another enclose an Eigenvalue solution at their core
- also the concept of the graph Laplacian reappears in spectral clustering, Laplacian Eigenmaps, "Manifold Forests"
- in both graphical models an underlying question is how dependent variables are modeled. In both cases we use the notion of Markovianity.

Notes on MRFs (paper on github: "ermongroup")  
and graph factorization

Clique: fully connected subgraph



$$p(x_i, x_j)$$



$$\Sigma(x_i)$$

$\Rightarrow$  MRFs with a maximum clique size of 2 + submodularity condition can be readily (efficiently) solved via graph cuts.

$\Rightarrow$  Cliques larger than 3 can in general not efficiently be solved