# Lecture 10

Hierarchical clustering
( NOT distances )

SVM - linear separable primal

- HW 2 - pb2 removed
Project t.b.p.
LV   SVM
LV.1   RKHS

Bootstrap
· Boosting
GP, RFF, DD
· Mani L ; DP mix

# Lecture IV – Hierarchical clustering. Comparing clusterings

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

STAT/BIOST 527

# Hierarchical Methods of Clustering

- **Agglomerative** (bottom up):
  - Initially, each point is a cluster
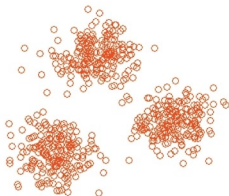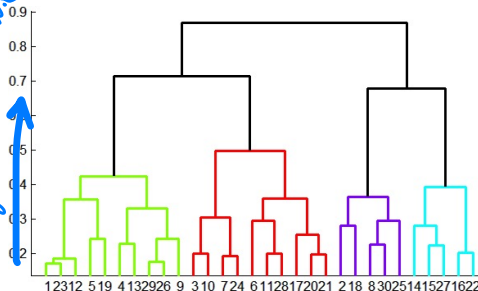  - Repeatedly combine the two "nearest" clusters into one

} greedy

- **Divisive** (top down):
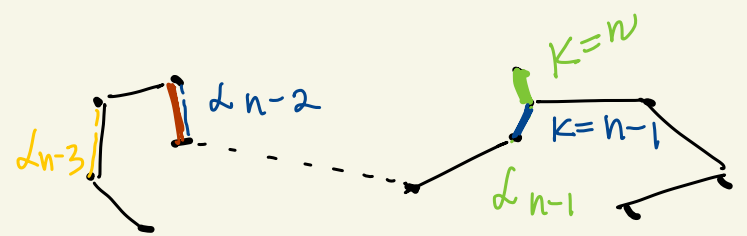  - Start with one cluster and recursively split it

Divisive

$\underline{\underline{Loss}}$ $\quad \alpha(\Delta_K) =$ at level $K$ $\qquad\qquad K = 1:n$

1. Single Linkage $\alpha_k \equiv \alpha(\Delta_k) = -\min\limits_{\substack{k \neq k'}} \min\limits_{\substack{i \in C_k \\ j \in C_{k'}}} \|x^i - x^j\|$

$\underline{\text{Min Spanning }\underline{\text{Tree}}}$

$\downarrow$     no cycles

contains all nodes
$\equiv$ connected
$\equiv n-1$ edges

log
sum
edge
length

Agglomeration

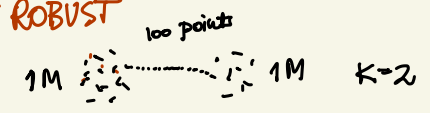$\alpha_{n-3}$   $\alpha_{n-2}$   $K = n$   $K = n-1$   $\alpha_{n-1}$

$K = n$
$K = n-1$
$K = n-2$
$\cdots$

$\Rightarrow$ tree (no cycles)
(MST)
Min Spanning Tree

$\Downarrow$

can be run as
Divisive too !!

Pb: NOT ROBUST
    100 points
1M $\cdots$ 1M   $K=2$

<u>Loss</u> $\alpha(\Delta_K)$ = at level $K$      $K = 1 : n$

1. Single Linkage $\alpha_k \equiv \alpha(\Delta_k) = -\min_{k \neq k'} \min_{\substack{i \in C_k \\ j \in C_{k'}}} \| x^i - x^j \|$

2. Least squares

$$\alpha(\Delta_K) = \sum_{k=1}^{K} \sum_{i \in C_k} \| x^i - \mu_k \|^2$$
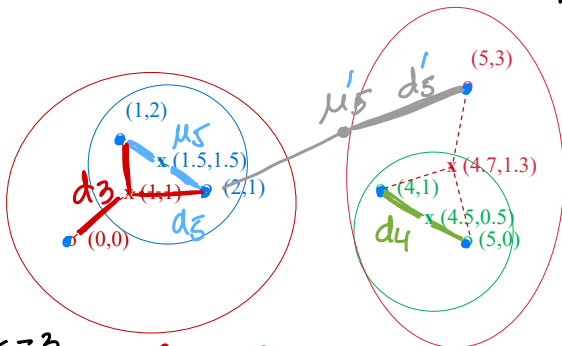
3. Mixture log-likelihood

$$\alpha(\Delta_k) = \sum_{k=1}^{K} \sum_{i \in C_k} \ln\left( \pi_k \, f_k(x^i) \right) \qquad \text{CONDITIONAL log-likelihood}$$

$$\ln \bar{u}_k + \ln f_k(x^i)$$

# Hierarchical clustering – Overview

(Dendrograms)

- **Agglomerative** (bottom up)
  - **Single linkage**
    - based on Minimum Spanning Tree
    - $\mathcal{O}(n^2 \log n)$
    - sensitive to outliers
  - Heuristics – average linkage
  - **Agglomerative K-means**
    - Loss $\mathcal{L}(\Delta_K) = 0$ for $K = n$
    - When $K \leftarrow K - 1$ (two clusters merged), $\mathcal{L}$ increases
    - For $K = n, n - 1, \ldots 2$, iteratively merge the 2 clusters that minimize increase of $\mathcal{L}$
    - $\mathcal{O}(n^3)$ – too expensive for big data

- **Divisive** (bottom down)
  - Recursively split $\mathcal{D}$ into $K = 2$ clusters
  - almost any clustering algorithm (e.g. K-means, min diameter)
  - notable example Spectral clustering (later)

  - Advantages
    - most important splits are first
    - can stop after only a few splits

# Example: Hierarchical clustering



$n = 6$

$K = 6$ $\mu_k = x^k$

$\alpha_6 = 0$

$k = 5$ $\alpha_5 = 0 \times 4 +$ $2 d_5^2$ BAD

$\alpha_5 = 0 \times 4 + 2 d_5^2$

$\alpha_4 = 0 \times 2 + 2 d_5^2 + 2 d_4^2$

$k = 3$

$\alpha_3 = 0 + 2 d_4^2 + 3 d_3^2 = \alpha_4 - 0 - 2 d_5^2 + 3 d_3^2$

Labels in figure: (5,3), (1,2), $\mu_5'$, $d_5'$, x (1.5,1.5), $\mu_5$, $d_3$ x (1,1), (2,1), $d_5$, (0,0), x (4.7,1.3), (4,1), x (4.5,0.5), $d_4$, (5,0)
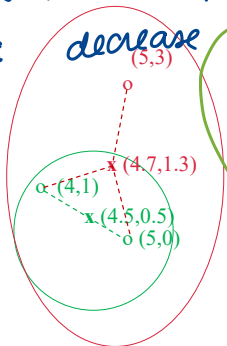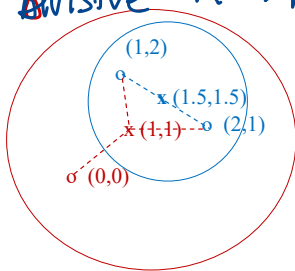
**Data:**
o ... data point
x ... centroid

**Dendrogram**

# Example: Hierarchical clustering

Agglomerative    $K \rightarrow K-1$    increase $d$ least

Divisive    $K \rightarrow K+1$    decrease $d$ most

(1,2)
o
x (1.5,1.5)
x (1,1)    o (2,1)
o (0,0)

(5,3)
o
x (4.7,1.3)
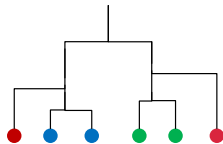o (4,1)
x (4.5,0.5)
o (5,0)

$K = 1 : 2^n$ options

$K = n : \dfrac{n(n-1)}{2}$ opti

**Data:**
o … data point
x … centroid

**Dendrogram**

# Lecture V: Support Vector Machines and Kernel Machines

*Isabelle Guyon*
*Corinna Cortes*
*Vladimir Vapnik*

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

1. max margin
2. convex opt
3. kernel trick

April, 2023

# Linear SVM's  ⬅

    The margin and the expected classification error    *why*

    Maximum Margin Linear classifiers  ⬅— 1,2

    Linear classifiers for non-linearly separable data

# Non linear SVM

    The "kernel trick"

    Kernels

    Prediction with SVM

# Extensions

    $L_1$ SVM

    Multi-class and One class SVM

    SV Regression

**Reading** AoNPS Ch.: Ch. 12.1–3, HTF Ch.: Ch 14 (14.1,14.2–14.2.4 kernels, 14.4 and equations
(14.28,14.29) kernel trick, 14.5.1.–3 Support Vector Machines)7.1–7.4, 7.7
Additional Reading: C. Burges - "A tutorial on SVM for pattern recognition"
These notes: Appendices (convex optimization) are optional.

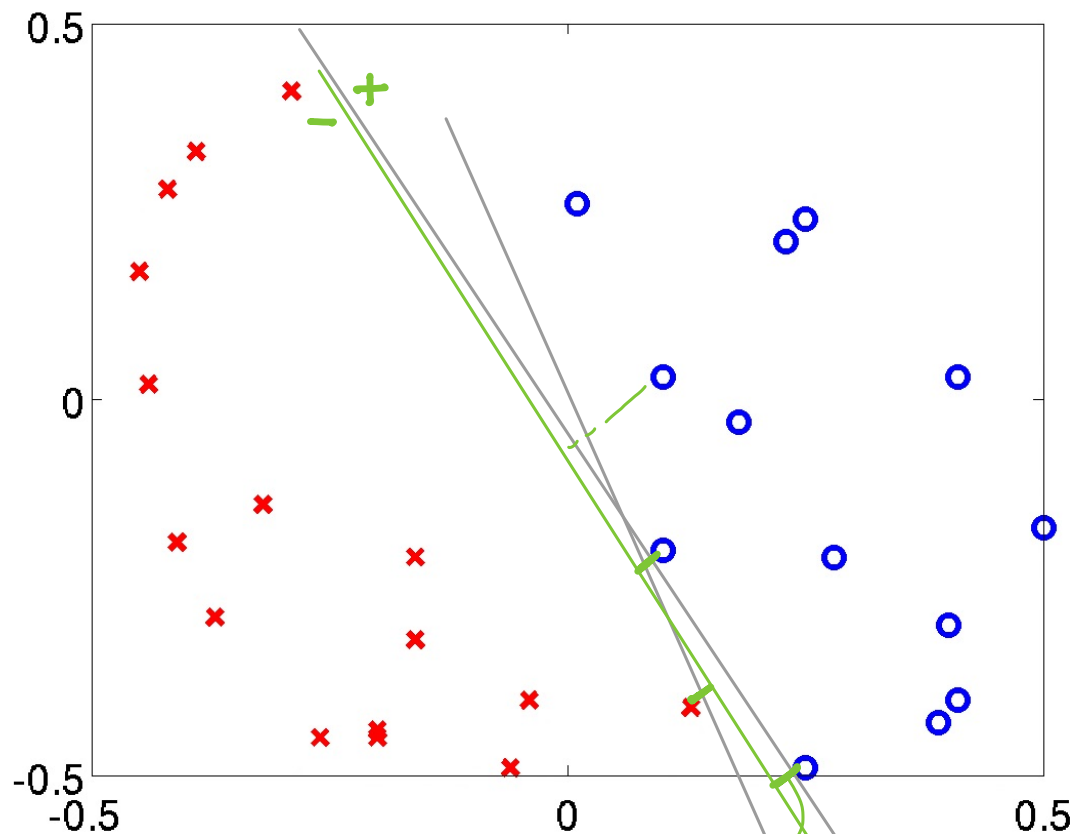Data Linearly separable ⇒ ∞ linear classifiers
How to choose $\hat{w}$ ?

$f(x) = w^T x + b$
$x, w \in \mathbb{R}^d$
$\hat{y} = \text{sign } f(x)$

Idea 1. Max Margin

Parametric
LDA
LR
Perceptron

$\rho$ = margin of $f_{w,b}$

− +

In $\mathbb{R}^d$ : $d+1$ points to determine $(W^*, b^*)$ = max margin hyperplane

             ↓

    support vectors

          "
          points in $\mathbb{R}^d$

SVM = max margin classifier

   "
  classifier

Robustness :
- any $x^i$ can be perturbed by $\leq \rho$ and not change label
- all $x^i$'s NOT SV - can be perturbed more
                          - have no influence on $W^*, b^*$
                                ⇒ fewer params ?
                         YES
- Theorem (s)

# The margin and the expected classification error

**Theorem** Let $\mathcal{F} = \{\operatorname{sgn}(w^T x), \|w\| \leq \Lambda, \|x\| \leq R\}$ and let $\rho > 0$ be any "margin". Then for any $f \in \mathcal{F}$, w.p $1 - \delta$ over training sets

*generalization err*

$$L_{01}(f) \leq \hat{L}_\rho + \sqrt{\frac{c}{n}\left(\frac{R^2\Lambda^2}{\rho^2}\ln n^2 + \ln\frac{1}{\delta}\right)} \qquad \frac{1}{\sqrt{n}}\sqrt{\frac{1}{\delta^2} + \cdots} \quad (5)$$

*classif err on $\mathcal{D}$*

where $c$ is a universal constant and $\hat{L}_\rho$ is the fraction of the training examples for which
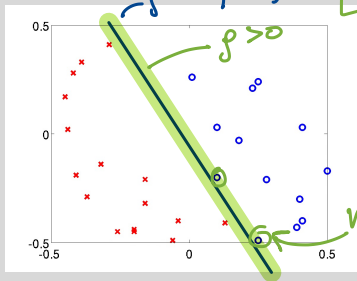
*including*
*margin errors*

$$y^i w^T x_i < \rho \qquad (6)$$

▶ a data point $i$ that satisfies (6) for some $\rho$ is called a **margin error**
▶ For $\rho = 0$ the margin error rate $\hat{L}_\rho$ is equal to $\hat{L}_{01}$

# The margin and the expected classification error

**Theorem** Let $\mathcal{F} = \{\text{sgn}(w^T x), \|w\| \leq \Lambda, \|x\| \leq R\}$ and let $\rho > 0$ be any "margin". Then for any $f \in \mathcal{F}$, w.p $1 - \delta$ over training sets

$$L_{01}(f) \leq \hat{L}_\rho + \sqrt{\frac{c}{n}\left(\frac{R^2\Lambda^2}{\rho^2}\ln n^2 + \ln\frac{1}{\delta}\right)} \qquad (5)$$

where $c$ is a universal constant and $\hat{L}_\rho$ is the fraction of the training examples for which

$$y^i w^T x_i < \rho \qquad (6)$$

▶ a data point $i$ that satisfies (6) for some $\rho$ is called a **margin error**
▶ For $\rho = 0$ the margin error rate $\hat{L}_\rho$ is equal to $\hat{L}_{01}$

Another theorem

$$\mathcal{F}_\rho = \{ f_{w,b} , \quad w, x \in \mathbb{R}, \quad \text{margin} \geq \rho \}$$

$$vc\,dim = min\{ d+1, \frac{1}{\rho^2} \}$$

# Maximum Margin Linear classifiers   *How to estimate?*

Support Vector Machines appeared from the convergence of Three Good Ideas
**Assume** (for the moment) that the data are linearly separable.

*all $x^i$ away from bdary*

▶ Then, there are an infinity of linear classifiers that have $\hat{L}_{01} = 0$. Which one to choose?

**First idea** Select the classifier that has **maximum margin** $\rho$ on the training set.

  ▶ For any parameters $(w, b)$ that perfectly classify the data $\hat{L}(w, b) = 0$.
  ▶ Among these, the best $(w, b)$ is the one that minimizes $\rho$ in 5
  ▶ Hence, we should choose

$$\underset{\rho, w, b: \hat{L}(w,b)=0}{\operatorname{argmax}} \quad \rho, \quad \text{s.t. } d(x, H_{w,b}) \geq \rho \text{ for } i = 1 : n, \tag{7}$$

*① wanted*   *all $(x^i, y^i)$ correct*

where $d()$ denotes the Euclidean distance and $H_{w,b} = \{ x \mid w^T x + b = 0 \}$ is the decision boundary of the linear classifier.

*↳ hyperplane*

▶ Because $d(x, H_{w,b}) = \frac{|w^T x + b|}{||w||}$ (proof in a few slides) (7) becomes

*②*
$$\underset{\rho, w, b: \hat{L}(w,b)=0}{\operatorname{argmax}} \quad \rho, \quad \text{s.t. } \frac{|w^T x^i + b|}{||w||} \geq \rho \text{ for } i = 1 : n, \tag{8}$$

*replace $d( )$*

# Maximum Margin Linear classifiers

sign $w^T x + b = y$     $y \in \pm 1$

$y(w^T x + b) > 0$

We continue to transform (8)

▶ If all data correctly classified, then $y^i(w^T x^i + b) = |w^T x^i + b|$. Therefore (8) has the same solution as

③

$$\operatorname*{argmax}_{\rho, w, b} \rho, \quad \text{s.t.} \quad \frac{y^i(w^T x^i + b)}{||w||} \geq \rho \text{ for } i = 1:n, \tag{9}$$

▶ Note now that the problem (9) is underdetermined. Setting $w \leftarrow Cw, b \leftarrow Cb$ with $C > 0$ does not change anything.

▶ We add a cleverly chosen constraint to remove the indeterminacy; this is $||w|| = 1/\rho$, which allows us to eliminate the variable $\rho$. We get

④

$$\operatorname*{argmax}_{w, b} \frac{1}{||w||}, \quad \text{s.t.} \quad y^i(w^T x^i + b) \geq 1 \text{ for } i = 1:n, \tag{10}$$

Note: the successive problems (7),(8),(9),... are equivalent in the sense that their optimal solution is the same.

$$\operatorname*{argmax}_{w, b} \frac{1}{||w||} \quad \text{s.t} \quad \frac{y^i(w^T x^i + b)}{||w||} \geq \frac{1}{||w||}$$

$\underbrace{\frac{1}{||w||}}_{\rho}$     $\underbrace{\frac{1}{||w||}}_{\rho}$

# Alternative derivation of (10)

**est idea** Select the classifier that has maximum margin on the training set, by the alternative definition of margin.

Formally, define $\min_{i=1:n} y^i f(x^i)$ be the **margin of classifier $f$ on $\mathcal{D}$**. Let $f(x) = w^T x + b$, and choose $w, b$ that

$$\text{maximize}_{w \in \mathbb{R}^n, b \in \mathbb{R}} \min_{i=1:n} y^i (w^T x^i + b) \; s.t. \; \hat{L}(w, b) = 0$$

▶ Remarks

  ▶ (if data is linearly separable), there exist classifiers with margins $> 0$
  ▶ one can arbitrarily increase the margin of such a classifier by multiplying $w$ and $b$ by a positive constant.
  ▶ Hence, we need to "normalize" the set of candidate classifiers by requiring instead

$$\text{maximize} \min_{i=1:n} d(x, H_{w,b}), \text{ s.t. } y^i (w^T x^i + b) \geq 1 \text{ for } i = 1 : n, \qquad (11)$$

  where $d()$ denotes the Euclidean distance and $H_{w,b} = \{ x \mid w^T x + b = 0 \}$ is the decision boundary of the linear classifier.
  ▶ Under the conditions of (11), because there are points for which $|w^T x + b| = 1$, maximizing $d(x, H_{w,b})$ over $w, b$ for such a point is the same as

$$\max_{w,b} \frac{1}{||w||}, \text{ s.t. } \min_i y_i (w^T x + b) = 1 \qquad (12)$$

# Second idea

(5) $\min\limits_{w,b} \|w\| \ \text{st} \ y^i(w^T x^i + b) \geq 1 \ \text{for} \ i = 1:n$

The **Second idea** is to formulate (10) as a **quadratic** optimization problem.

(6) $$\min\limits_{w,b} \frac{1}{2}\|w\|^2 \ \text{s.t} \ y^i(w^T x^i + b) \geq 1 \ \text{for all } i = 1:n \tag{13}$$

This is the **Linear SVM (primal) optimization problem**

▶ This problem has a strongly convex **objective** $\|w\|^2$, and **constraints** $y^i(w^T x^i + b)$ linear in $(w, b)$.

▶ Hence this is a convex problem, and can be studied with the tools of convex optimization.

objective
quadratic
convex

constraints
linear
$\Rightarrow \{$ feasible $w,b\}$ convex

# The distance of a point $x$ to a hyperplane $H_{w,b}$

$$d(x, H_{w,b}) = \frac{|w^T x + b|}{||w||} \tag{14}$$

Intuition: denote

$$\tilde{w} = \frac{w}{||w||}, \ \tilde{b} = \frac{b}{||w||}, \ x' = \tilde{w}^T x. \tag{15}$$

Obviously $H_{w,b} = H_{\tilde{w},\tilde{b}}$, and $x'$ is the length of the projection of point $x$ on the direction of $w$.

The distance is measured along the normal through $x$ to $H$; note that if $x' = -\tilde{b}$ then $x \in H_{w,b}$ and $d(x, H_{w,b}) = 0$; in general, the distance along this line will be $|x' - (-\tilde{b})|$.