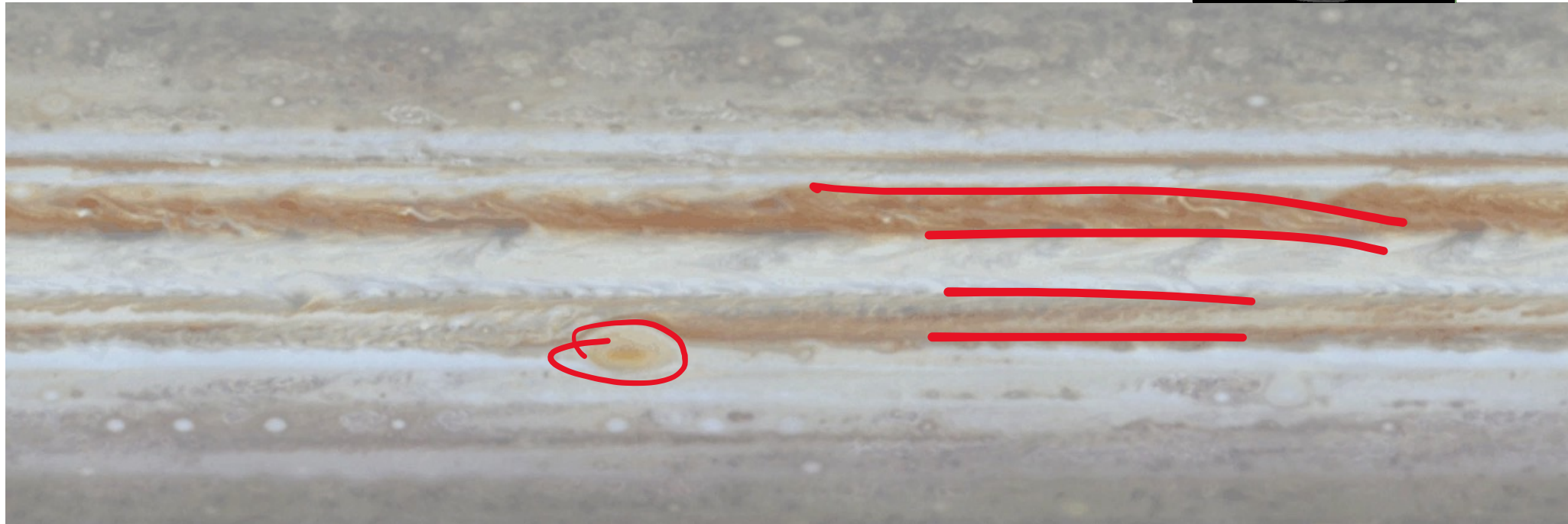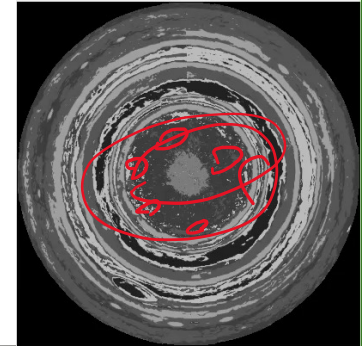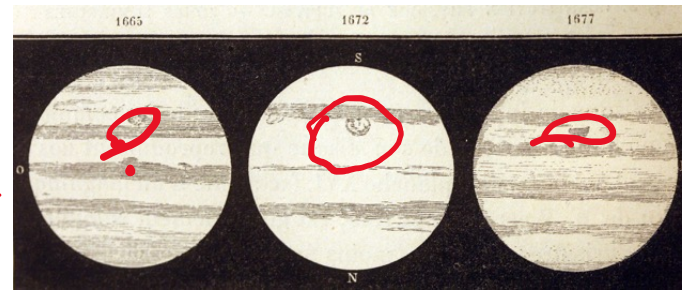# EE520 Data Driven Analysis of Complex Systems
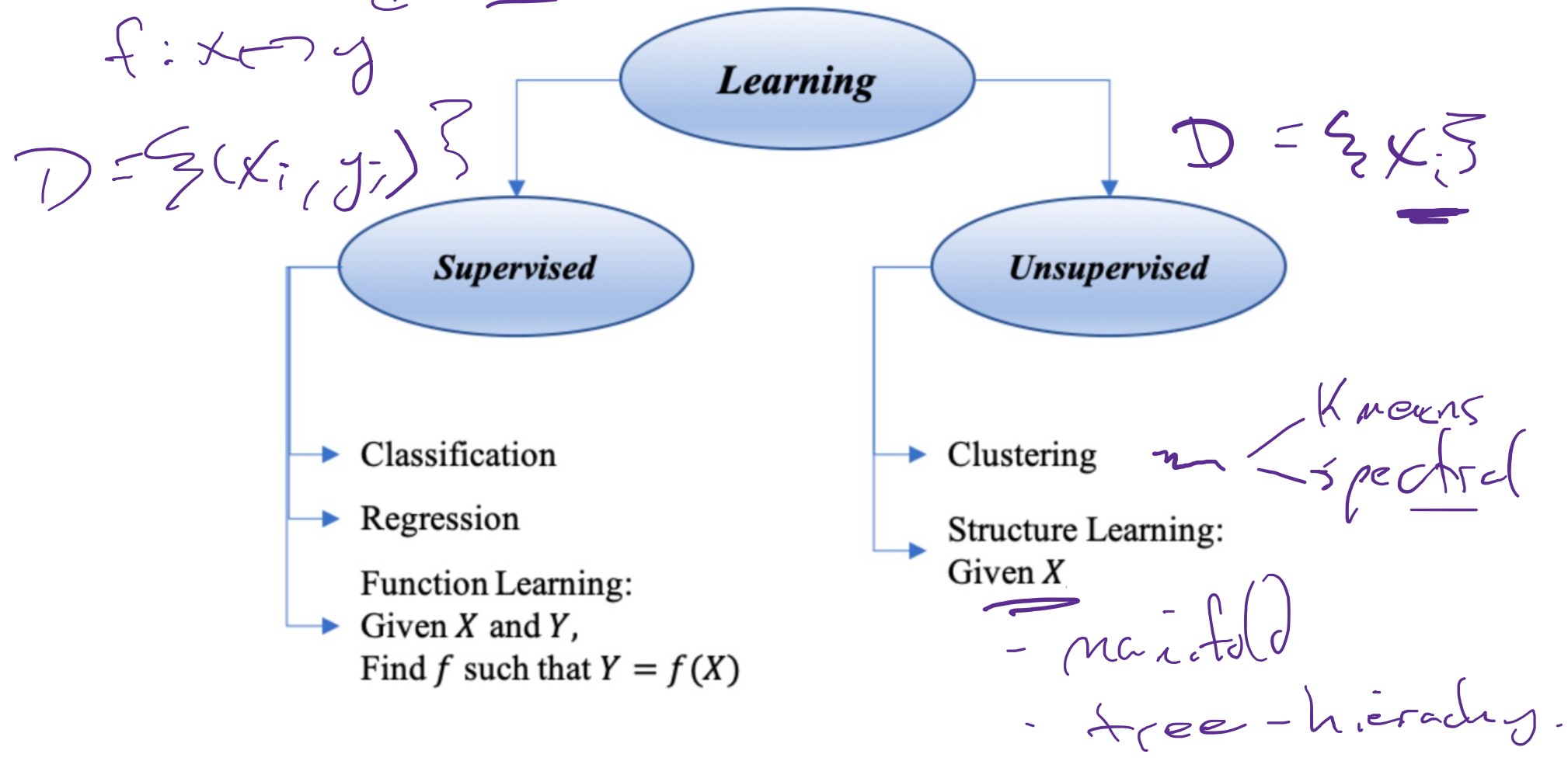
Erik Bollt

# Ch 5 - Clustering and Classification

attributes

$f: x \rightarrow y$ labels.

$D = \{(x_i, y_i)\}$



**Learning**

**Supervised**

**Unsupervised**

Classification

Regression

Function Learning:
Given $X$ and $Y$,
Find $f$ such that $Y = f(X)$

Clustering
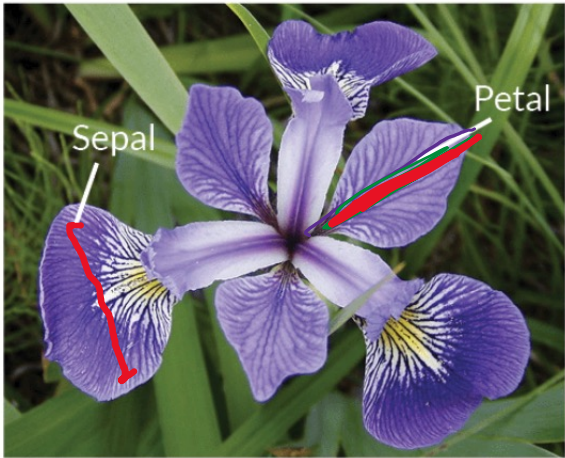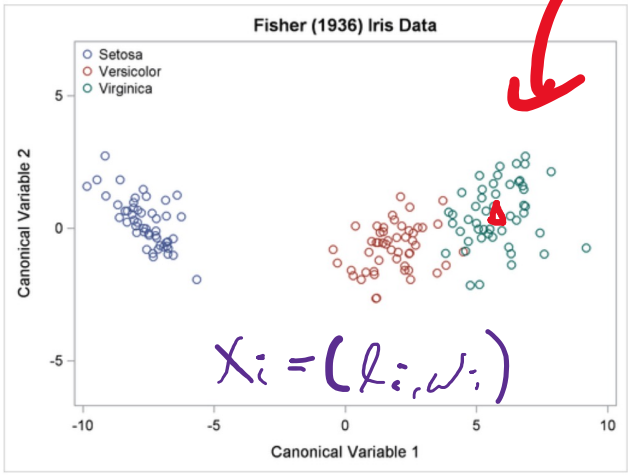
Structure Learning:
Given $X$

$D = \{x_i\}$

K means
spectral

- manifold
- tree-hierarchy.

Discriminating Fisher's iris data by using the petal areas

The Iris Dataset contains four features (length and width of sepals and petals) of 50 samples of three species of Iris (Iris setosa, Iris virginica and Iris versicolor). These measures were used to create a linear discriminant model to classify the species. The dataset is often used in data mining, classification and clustering examples and to test algorithms.
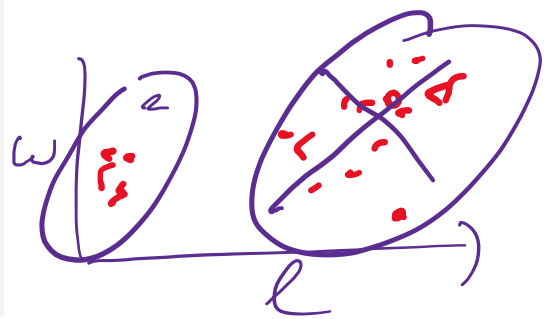


$$X_i = (\ell_i, w_i)$$
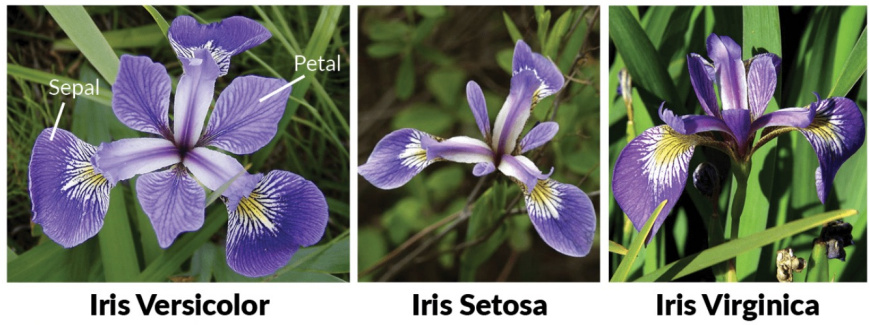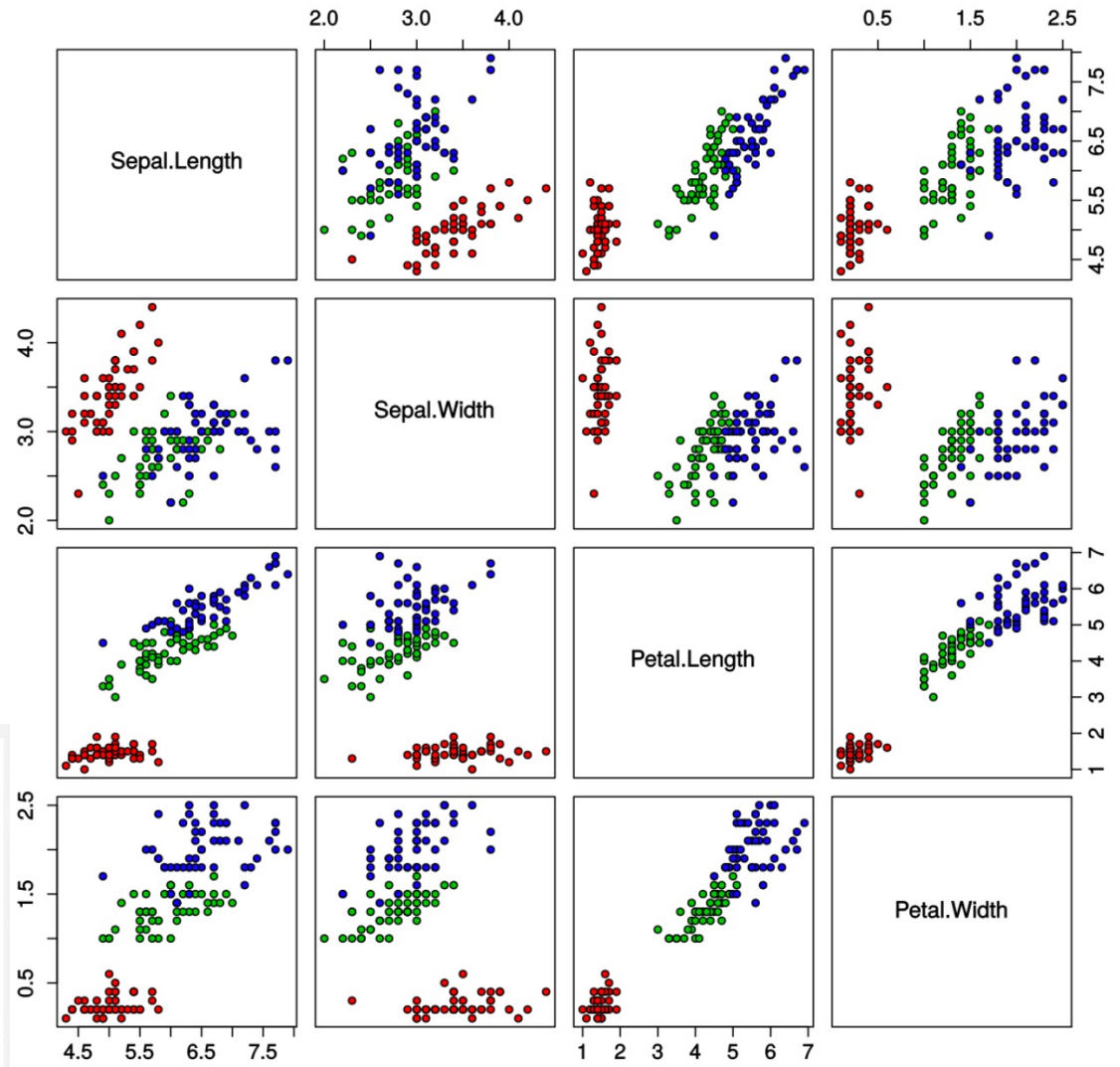
Ronald Fischer 1936



Iris Data (red=setosa,green=versicolor,blue=virginica)



**Iris Versicolor**     **Iris Setosa**     **Iris Virginica**

**Data** **Learning** $D_1$

$D_2$

$x_i = \{$ engine, zipcode, time $\delta f \}$

$y_i = \{0, 1, 2, 3\}$

$f: x \to y$

**supervised** learn $f$ **unsupervised** $\perp$ no $y_i$, no $f$

$\llcorner$ regression $(x_i, y_i) \in \mathbb{R}^q \times \mathbb{R}^d$ $\hookrightarrow$ just structure

$\llcorner$ classify if $y_i \in Z_r = \{0, 1\}$ just shape $= \{$ hot, cold $\}$

$f: \underline{X} \to \underline{Y}$

$D_2 = \{(x_i, y_i)\}_i$ $x_i \sim \underline{X}$, $y_i \sim \underline{Y}$

either vs. $z_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix}$

$D_1 = \{x_i\}_{i=1}^N$, $x_i \sim \underline{X}$

If we are going to do some machine learning – we had better get serious about what does learning mean?

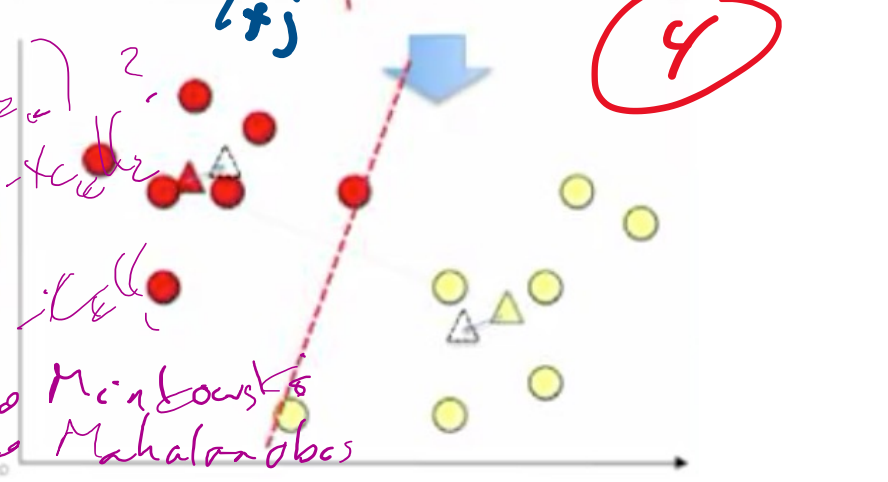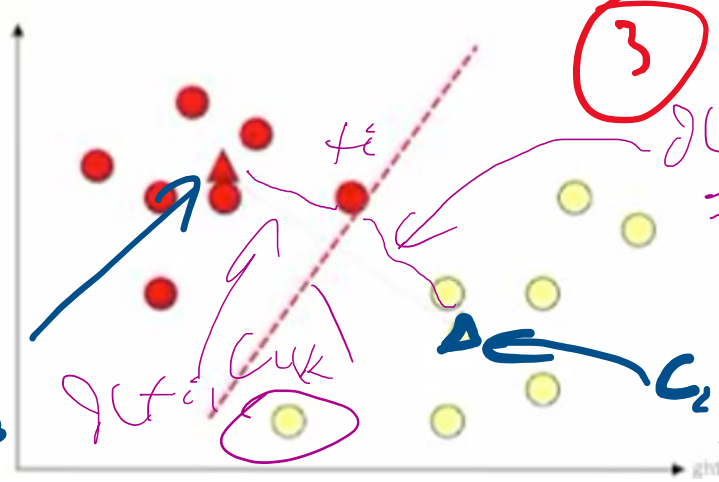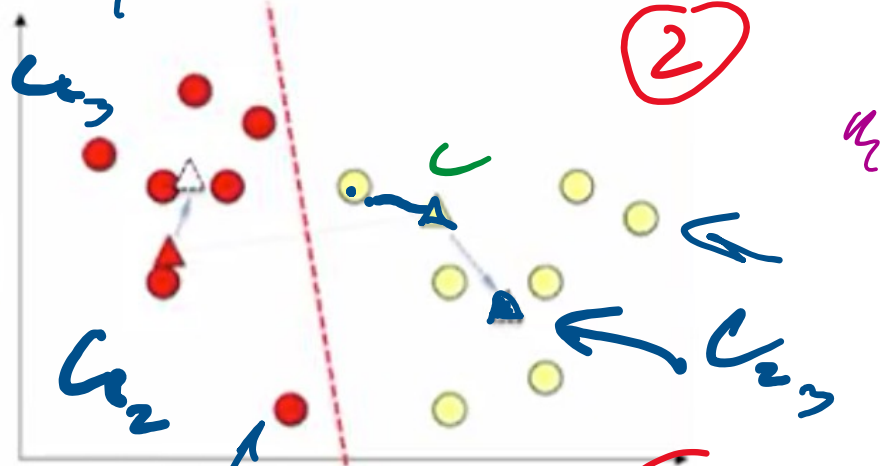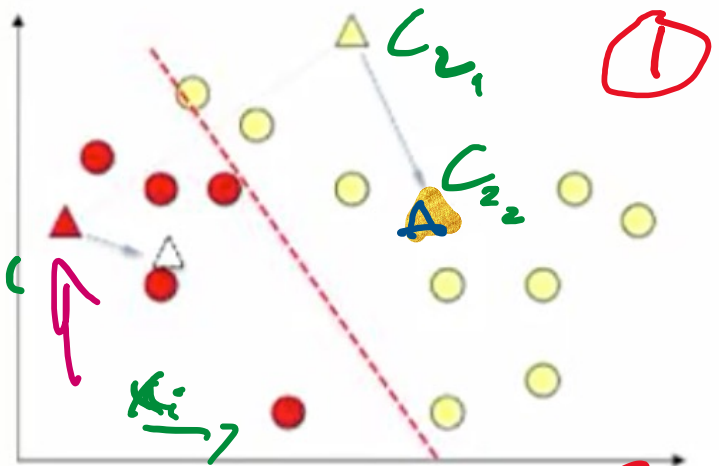-Supervised discrete output (labels)
-supervised continuous output
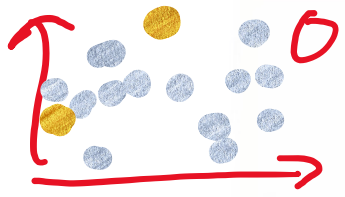
-unsupervised discrete output (labels – cluster analysis)
-unsupervised continuous output (ROM, manifold learning, density estimation)

Learning functions?
Learning structure?

# K-means clustering example



$C_{2_1}$ $C_{2_2}$

① ② ③ ④

$C_{2_3}$

$x_i$

$C_{1_3}$

$d(x_i, C_{2_K})$

$d(x_i, C_{2_K})^2 = ||x_i - x_{C_2}||_2$

$||x_i - x_j||_2^2$

Minkowski

Mahalanobis

https://files.realpython.com/media/centroids_iterations.247379590275.gif
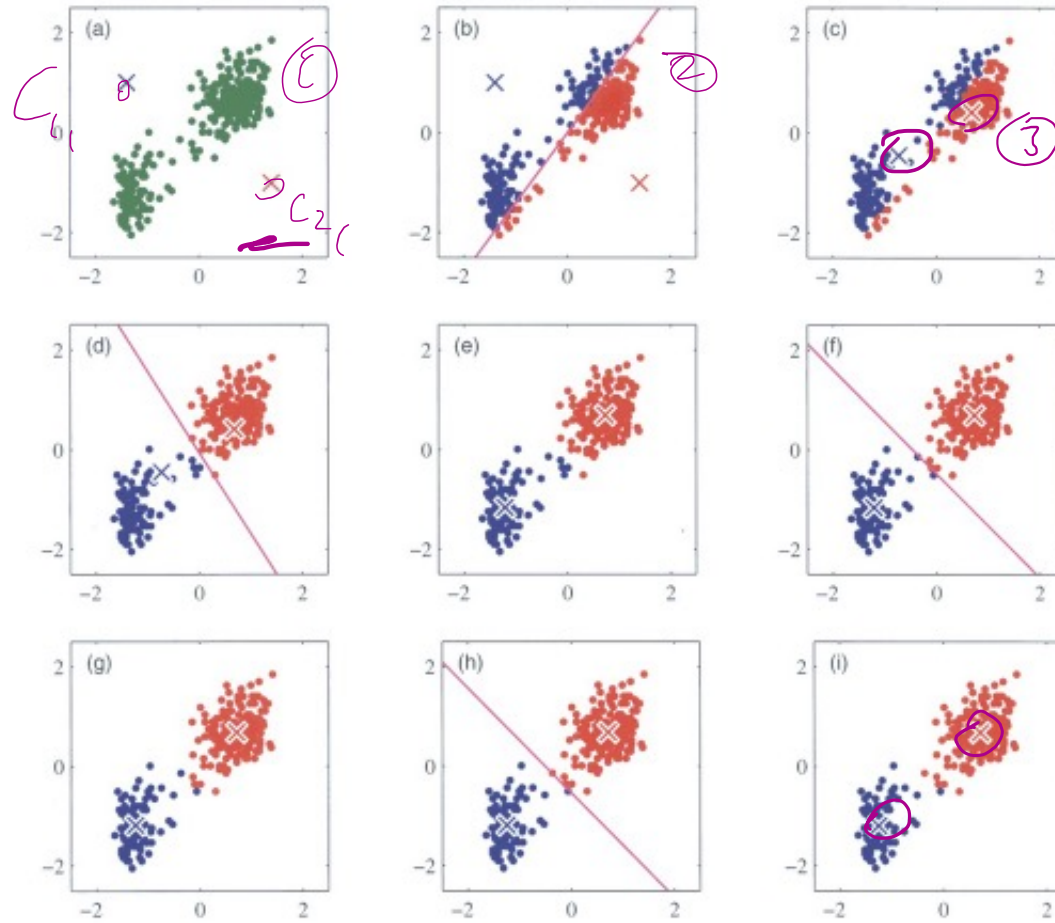
$x_i = (x_{i_1}, x_{i_2})$



Illustration of K-means algorithm from [Bishop 2006]

# Kmeans Convergence

**Objective**

$$\min_{\mu}\min_{C} \sum_{i=1}^{k} \sum_{x \in C_i} |x - \mu_i|^2$$

1. Fix $\mu$, optimize $C$:

$$\min_{C} \sum_{i=1}^{k} \sum_{x \in C_i} |x - \mu_i|^2 = \min_{C} \sum_{i}^{n} |x_i - \mu_{x_i}|^2$$

**Step 1 of kmeans**

2. Fix $C$, optimize $\mu$:

$$\min_{\mu} \sum_{i=1}^{k} \sum_{x \in C_i} |x - \mu_i|^2$$

 – Take partial derivative of $\mu_i$ and set to zero, we have

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

**Step 2 of kmeans**

Kmeans takes an alternating optimization approach, each step is guaranteed to decrease the objective – thus guaranteed to converge

[Slide from Alan Fern]

# K-Means

- An iterative clustering algorithm

  - Initialize: Pick $K$ random points as cluster centers

  - Alternate:
    1. Assign data points to closest cluster center
    2. Change the cluster center to the average of its assigned points

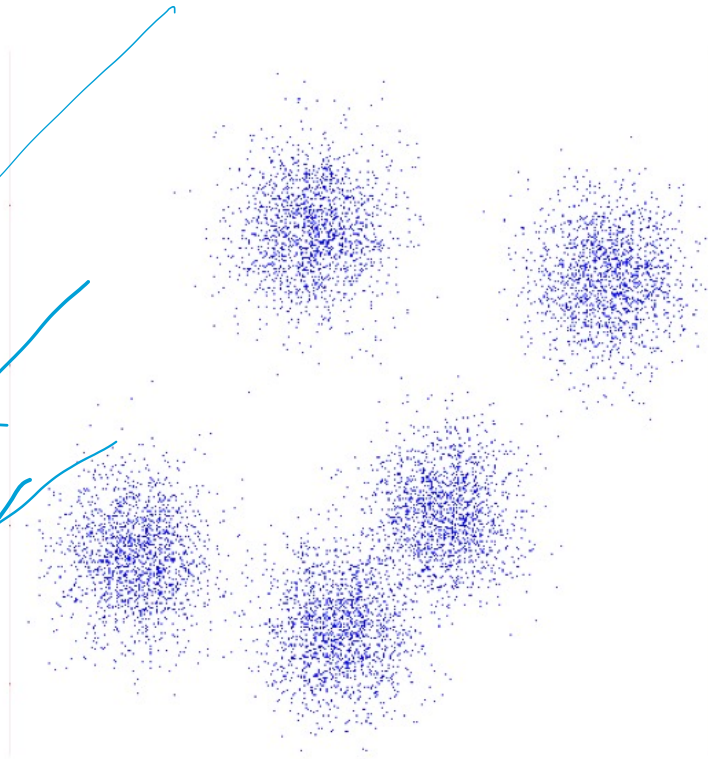  - Stop when no points' assignments change

# K-Means

- An iterative clustering algorithm

  - Initialize: Pick *K* random points as cluster centers
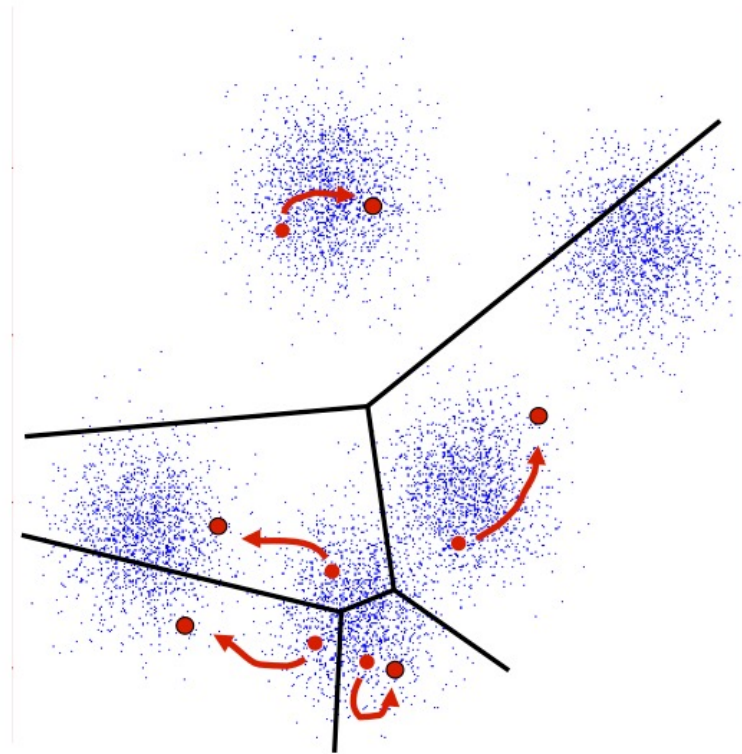
  - Alternate:
    1. Assign data points to closest cluster center
    2. Change the cluster center to the average of its assigned points
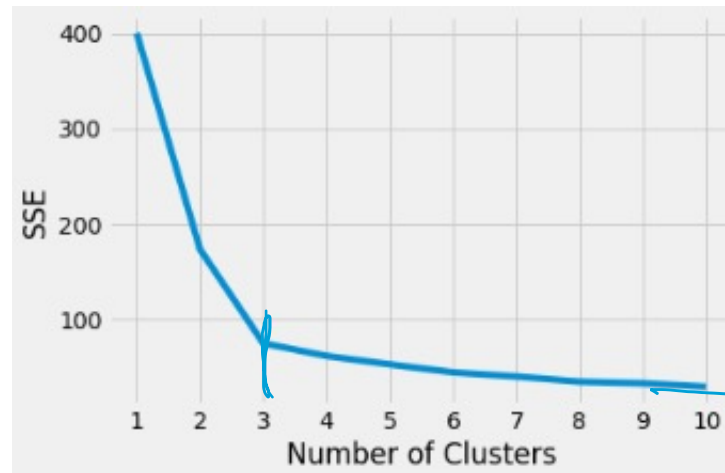
  - Stop when no points' assignments change

elbow

$$X = (x_1, x_2)$$
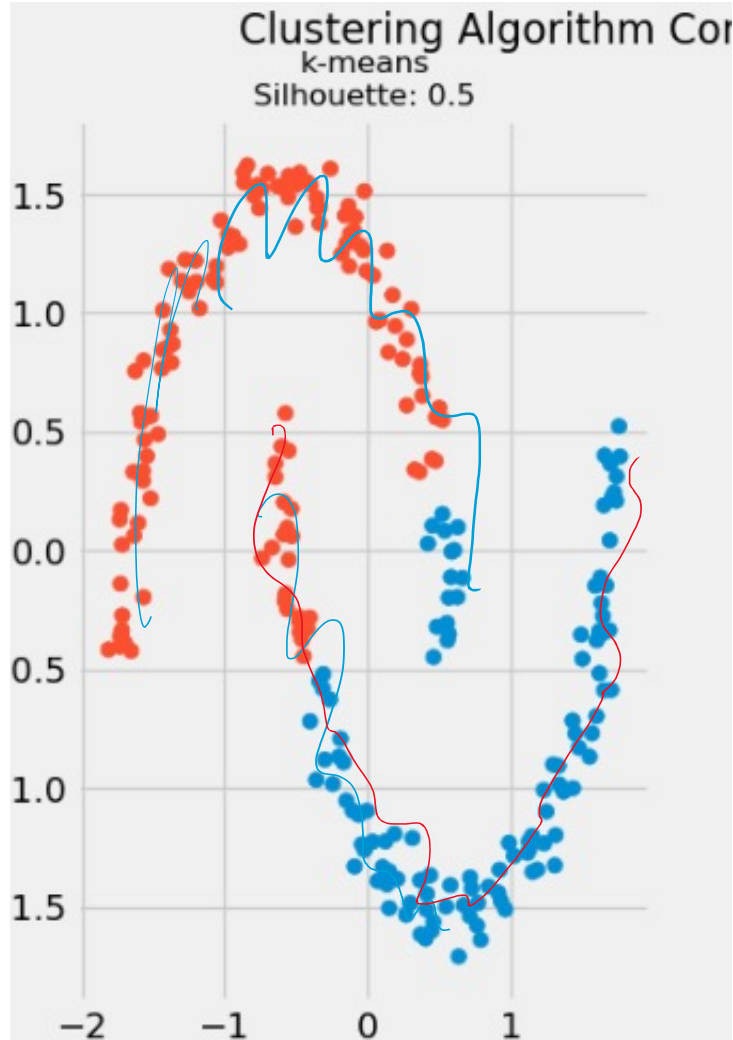
$$C = (C_1, C_2)$$

$$d(x, c)^2$$

$$= \sum_i (x_i - c_i)^2$$

$$= \frac{\sum (x_i - c_i)^4}{\sum e^{*(x_i - c_i)}} \; ?$$



Clustering Algorithm Co[...]
k-means
Silhouette: 0.5

# K-means not able to properly cluster



$x_i = (2, -3)$

$\rightarrow (\theta, r)$

$\rightarrow (300°, 2.5)$

# Changing the features (distance function) can help

# Example: K-Means for Segmentation



K=2    K=3    K=10    Original

# Image segmentation

Goal: Break up the image into meaningful or perceptually similar regions

# Clustering gene expression data

Eisen et al, PNAS 1998

# Cluster news articles

# Cluster people by space and time



[Image from Pilho Kim]

$$z^i = ((x_1^i, y_1^i, t_1^i), (x_2^i, y_2^i, t_2^i), (x_3^i, y_3^i, t_3^i))$$

$$\{z_i\}_{i=1}^{N}$$

Clustering on tracking moving targets – prepare a vector – feature - corresponding to position over time.



$$X_i(t_N)$$

$$\left\{ X_i(t_j) \right\}_{\substack{j=1..N \\ i=1..M}}$$

$$X_i(t_2) \quad \underline{X_i(t_3)}$$

$$\underline{X_i(t_1)}$$

$$\underset{2N\times1}{\longrightarrow} \underline{E_i} = [ \overset{2\times1}{x_i(t_1)} \ldots x_i(t_N) ]$$

# Clustering species ("phylogeny")



[Lindblad-Toh et al., Nature 2005]

# Brain Tumor Detection and Identification Using K-Means Clustering Technique

**Malathi R**
Department of Computer Science, SAAS College, Ramanathapuram, Email: malapraba@gmail.com
**Dr. Nadirabanu Kamal A R**

Original Image    Segmented Image with #partitions = 5    Extracted Image with label = 2

Medical imaging – cancer tissue

# DETECTING AND COUNTING THE NO. OF WHITE BLOOD CELLS IN BLOOD SAMPLE IMAGES BY COLOR BASED K-MEANS CLUSTERING

[1]Neha Sharma, [2]Nishant Kinra
[1]Indira Gandhi Delhi Technical University for Women, New Delhi, India
[2]Deenbandhu Chhotu Ram University of Science and Technology, Haryana, India
[1]neha.sksharma@yahoo.co.in, [2]nishant.kinra@gmail.com



Fig (1): Example of Leukocytosis



Fig (2): Comparison between Normal Blood and Leukemia

So that's it for the time being with unsupervised learning, a few clustering methods – mostly kmeans,

But in your book are other interesting clustering methods in 5.4, 5.5, including tree methods and also Gaussian Mixture models that are very popular.

Later we will also develop a spectral clustering method that is very general and powerful.

Also in unsupervised learning later we will do "manifold learning".

But for now…..

We transition to supervised learning

$1^{st}$ in 5.6 we will cover Fischer's Linear Discriminant (LDA)

$2^{nd}$ in 5.7 on to support vector machines (SVM)

Then Chapter 6 a grandly popular method – artificial neural nets.

Supervised Learning – Cat or Dog?

$y_i = Dog = 1$

$x_i = i$th picture

$[x_i]_{qp} x_i =$

160



$y_i = cat$

100

200

$D = \{(x_i, y_i)\}_{i=1}^{N}$

$x_i = i$th picture

$y_i = i$th label

$f: X \mapsto Y$ ; $f: \mathbb{R}^{20,000} \to \mathbb{Z}_2$

To distinguish cat's from dogs – first it will be more efficient if **we choose good (efficient) features.**

(a)

(b)

(c)

(d)

Wavelets then PCA (SVD) will work well.

onsupervised.

$D = \{(x_i)\}$

structure

cluster

ROM


(a)


(b)


(c)


(d)

$D = \{(x_i, y_i)\}$ supervised — classify, $y = -1$ or $1$.

# Fischer's Linear Discriminant Analysis LDA



(a)

(b)

(c)

(d)

PCA$_4$

PCA$_2$

cats

dogs

poor discrimination

optimal projection

$$\mathbf{w} = \arg\max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

where the scatter matrices for between-class $\mathbf{S}_B$ and within-class $\mathbf{S}_W$ data are given by

$$\mathbf{S}_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$$

$$\mathbf{S}_W = \sum_{j=1}^{2} \sum_{\mathbf{x}} (\mathbf{x} - \mu_j)(\mathbf{x} - \mu_j)^T.$$

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w}$$

Rayleigh quotient whose solution can be found via the generalized eigenvalue problem

# Haar Wavelet

Define

$$\psi(x) \equiv \begin{cases} 1 & 0 \le x < \frac{1}{2} \\ -1 & \frac{1}{2} < x \le 1 \\ 0 & \text{otherwise} \end{cases}$$

and

$$\psi_{j\,k}(x) \equiv \psi\left(2^j x - k\right)$$

for $j$ a nonnegative integer and $0 \le k \le 2^j - 1$.

$\psi_{0,0} = \psi(x)$

$\psi_{1,0} = \psi(2x)$   $\psi_{1,1} = \psi(2x - 1)$

$\psi_{2,0} = \psi(4x)$   $\psi_{2,1} = \psi(4x - 1)$   $\psi_{2,2} = \psi(4x - 2)$   $\psi_{2,3} = \psi(4x - 3)$

Mother function.

dyadic

$D_k$

$(x_i, \psi_k)$

Mathworld
~ mathematica

$x_i(s) = s^2$

$x_i(s) \, \psi(s)$

$\int x_i(s)\, \psi_k(s)\, ds$

$\cos x$

$\cos 2x$

$\phi_1(s) = \psi_{00}$

$\phi_2(s) = \psi_{10}$

$\phi_3(s) = \psi_{1,1}$

# Better features

$\{\phi_k\}_{k=1}^{K} = \left\{ \begin{array}{l} \{\chi_k\} \quad \text{indicator functions.} \\ e^{iks} \quad \text{fourier basis} \leftarrow \\ \psi_k(s) \quad \text{haar wavelet basis} \end{array} \right.$

$\phi_k$

$e^{iks} = \cos ks + i \sin ks$

$a_k = \langle x_k, \phi_k \rangle = \int x_k(s) \phi_k(s) ds = \sum x_k(s) \phi_k(s)$

ith

Project onto basis function.

$\in \mathbb{R}^{20,000}$

$e_3$

$e_i$

$\phi_2$

$\phi_i$

$a_1$

$a_2$

$e_1$

$\phi_1$

Collection of pics as
in

$= \dfrac{20,000 \text{ the}}{\text{values in the pixels}}$

each pixel is considered
a feature.

Characteristic
fn: $\chi_k(s) = \begin{cases} 1 & s \in k \\ 0 \end{cases}$

$P_k = (x_i, \chi_k) = \int x_i(s) \chi_k(s)\,ds$

↑ feature

$\vec{P} = \begin{bmatrix} P_1 \\ \vdots \\ P_N \end{bmatrix} = $

$\in H = \mathbb{R}^{20,000}$  Data cloud.

$|C_k|^2 |a_k|^2$

Red-Fourier

"Power spectured."

$\omega$   $2\omega$   $4\omega$

$k$

$\int |X(s)|^2 \, ds \overset{\leq}{\approx} \sum |a_k|^2$  ⟵ Parseval

$b_1 = 0$
$b_3 = 0$
$a_j = (x :$ $\phi_j$
$b_{75} = -1$

$\vec{a}_i =$   $a_{1,i}$
$a_{2,i}$
$a_{3,i}$
$\vdots$
$a_{k,i}$

$$J(\omega) = \frac{\omega^T S_B \omega}{\omega^T S_W \omega} \quad ; \quad \omega = \arg\max_{\omega} J(\omega)$$

$$\min -\frac{1}{2}\omega^T S_B \omega \quad \longleftarrow$$

$$s.t. \quad \omega^T S_W \omega = 1$$

$$I \Rightarrow \tau$$

$$\phi = \frac{1}{2}\omega^T S_B \omega + \frac{\lambda}{2}(\omega^T S_W \omega - 1)$$

$$(\mu_2 - \mu_1) \cdot v$$
$$= (\mu_1 - \mu_2)^T v$$
$$= v \cdot (\mu_2 - \mu_1)$$
$$= v^T (\mu_2 - \mu_1)$$

$$\left((\mu_2 - \mu_1) \cdot w\right)^2 = w^T (\mu_2 - \mu_1) \cdot (\mu_2 - \mu_1) w$$

let
$$\mathcal{L} = -\frac{1}{2} \omega^T S_B \omega + \frac{1}{2} \lambda (\omega^T S_\omega \omega - 1) \Longleftarrow$$

$$\frac{\partial \mathcal{L}}{\partial \omega_1} = 0 \; ; \; \frac{\partial \mathcal{L}}{\partial \omega_2} = 0, \cdots \; ; \; \nabla_\omega \mathcal{L} = 0$$

$$\boxed{KKT}$$

$$\nabla_\omega \mathcal{L} = -S_B \omega + \lambda S_\omega \omega = 0$$

$$\boxed{\prod \quad S_B \underline{\omega} = \lambda S_\omega \underline{\omega}}$$

$$\boxed{\nabla \left( \frac{1}{2} \underline{\omega}^T \underline{\omega} \right) = \underline{\omega}}$$

$$\boxed{A x = \lambda B x} \quad \text{generalized eigenvector/value statement!}$$

$\omega$

$\omega$

$x$

$f : \overline{X} \rightarrow \overline{Y}$

$f : \overline{X} \rightarrow ?$

$(f, \omega)$

$$\underline{KKT}: \qquad x^\# = \underset{x}{\arg\min}\ f(x)$$

$$\text{subj.} \quad h_i(x) = 0, \quad \forall\ i=1,\dots,m\ ;\quad \vec{h}(\vec{x}) = \vec{0}$$

$$g_i(x) \leq 0, \quad \forall\ i=1,\dots,n \qquad \cancel{\vec{g}(\vec{x}) \leq \vec{0}}$$

$$x^\# = \underset{x,\lambda,\mu}{\arg\min}\ \mathcal{L}(x,\lambda,\mu) = \underset{x}{\arg\min}\ f(x) + \sum_i^m \lambda_i h_i(x) +$$

$$\sum_i^n \mu_j g_j(x)$$



$$f: \overline{X} \rightarrow \mathbb{R}$$

$\mathcal{I}$

Case A

not active
constraint

- or -

Case B

tangent.
$g_i(x) = 0$ ←
and f smallest.

$g_i(x) = 0$ .

Kkt, look for where level sets are tangent.

$g_i(x) = 0$

$-\nabla g$

Active

$g_i(x) \leq 0$ not active constraint

for

$\min f(x)$
subj

$g_i(x) \leq 0$ ←

$f$

$\nabla f$

$f$

KKT $\nabla g \parallel \nabla f$ at constrained opt.

$$\lambda \nabla g = \nabla f$$

$$\nabla f + \lambda \nabla g = 0 \quad \cdots \quad \text{Lagrange}.$$
$$\text{mult.}$$

Equality constr.

$$\rightarrow \nabla_\lambda f + \sum \nabla_\lambda \lambda_i h_i(x) + \sum \mu_i \nabla_\lambda g_i(x) = 0$$

- stationarity $\leftarrow$
- ineq. constr true.

# Support Vector Machines (SVM) — 5.7 Linear

Then
Nonlinear (kernelized) SVM (KSVM)

Wide Margin Decision Hyperplane for
Supervised - Learning Classification

Smola - Schökopff.

First a  linear binary classification – decision boundary/hyperplane

- Instance space: $x \in X$ ($|X| = n$ data points)
  - Binary or real-valued feature vector $x$ of word occurrences
  - $d$ features (words + other things, d~100,000+)
- Class: $y \in Y$
  - $y$: Spam (+1), Ham (-1)

$\{(x_i, y_i)\}_{i=1}^{n}$

$y_i = \frac{\pm}{2}$

$\{= 2$

| Viagra | Learning | The | Dating | Nigeria | Is_spam |
|--------|----------|-----|--------|---------|---------|
| 1 | 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 0 | 0 | -1 |
| 0 | 0 | 0 | 0 | 1 | 1 |

Review ... a hyperplane defined by a vector.

$\vec{n} = (a, b, c)$

$\vec{x} = (x_1, x_2, x_3)$

$\vec{x_0} = (x_{1,0}, x_{2,0}, x_{3,0})$

o l eqⁿ is a co-dim-1 restriction of space

$\vec{v} = \vec{x} - \vec{x_0} = \langle (x_1 - x_{10}), (x_2 - x_{20}), (x_3 - x_{30}) \rangle$

$\Pi := \{ x = (x_1, x_2, x_3) : \vec{v} = \vec{x} - \vec{x_0} \perp \vec{n} \}$

$v \cdot n = 0 \iff v \perp \vec{n}$

$n \cdot v = \langle a, b, c \rangle \cdot \langle x_1 - x_{10}, x_2 - x_{20}, x_3 - x_{30} \rangle = a(x_1 - x_{10}) + b(x_2 - x_{20}) + c(x_3 - x_{30}) = 0$

$SVM:$ $\quad D = \{(x_i, y_i)\}_{i=1}^{n}$ $\quad x_i \in \mathbb{R}^d$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad y_i \in \mathbb{Z}_2 = \{-1, 1\}$



$x_{i,2}$

Feature space

wide margin as big possible

$x_{i,1}$

hyper-plane

$x_2$ or $z$

$\vec{w}$

$-b\hat{j}+$

$x_i$
$y_i = 1$

- find $\vec{w}$ & $b$
- to optimize
  "wide margin"

$x_j$
$y_j = -1$

$\vec{w} \cdot x_i + b > 0$

$\vec{w} \cdot x_j + b < 0$

$b$

$x_1$

$\vec{w}$

$\vec{w} \cdot x + b = 0$

$\vec{w} \cdot x = 0$

$$(\vec{x}_i, y_i) \quad, \quad y_i = -1, \boxed{1}$$

$$\begin{pmatrix} y_i = 0 \text{ or } 1 \\ \text{apple or orange} \\ \text{dog or cat.} \end{pmatrix}$$

$$y_j(\omega \cdot x_j + b) = \text{sgn}(\omega \cdot x_j + b) \quad \text{good label.}$$

$$\text{sgn}(s) = \begin{cases} 1 & \text{if } s > 0 \\ 0 & \text{if } s = 0 \\ -1 & \text{if } s < 0 \end{cases}$$



$$\text{sgn}(\omega \cdot x_j + b) = 1 \quad \text{labelled well} $$
$$= -1 \quad \text{mislabelled.}$$

a loss function.

$$\ell(y_i, \bar{y}_i) = \ell(y_i, \text{sgn}(w \cdot \vec{x}_i + b)) = \begin{cases} 0 & \text{if correct label} \\ & y_i = \text{sgn}(w \cdot x_i + b) \\ 1 & \text{incorrect label} \end{cases}$$

$\bar{y}_i$

label you infer from $\vec{x}_i$ alone
if you have a good hyperplane $\vec{w}$ & $b$

Total loss

$$\sum_{i=1}^{N} \ell(y_i, \bar{y}_i)$$

Cost :  ... $\arg\min \left(\dfrac{1}{2} \|\omega\|_2^2\right)$ $\xrightarrow{\text{MSE}}$

subj $y_i(\omega^T x_i - b) - 1 = 0$

$\omega \cdot x_j$

small $\|\omega\|$

subj every matches truth

dist between   dist  big

$\dfrac{2}{\|\omega\|_2^2}$ big   $\approx$   $\boxed{\dfrac{\|\omega\|_2^2}{2}}$ small

constrained opt. $(\vec{w}, \underline{b})$ $\quad \vec{w} \in \mathbb{R}^d, d=2$ $\quad d+1 = 3$

$$\mathcal{L}(X, S, \Theta) = \frac{1}{2}\|w\|_2^2 - \sum_{i=1}^{l} S_i\left(y_i(w^T x_i - b) - 1\right)$$

$\{(x_i, y_i)\}$

$\frac{1}{2} w \cdot w$

$$y_1(w^T x_1 - b - 1) = 0$$
$$y_2(w^T x_2 - b) - 1 = 0$$
$$\vdots$$
$$y_n(w^T x_n - b - 1) = 0$$

$$= \frac{1}{2}\|w\|_2^2 - \sum_{i=1}^{l} S_i y_i (w^T \vec{x}_i - b) \quad \leftarrow \sum_{i=1}^{l} S_i$$

$$\nabla_\Theta \mathcal{L} = \vec{0}$$

$$\nabla_w \mathcal{L}(X, S, \Theta) = w - \sum S_i y_i \vec{x}_i = \vec{0} \quad \leftarrow d\text{-eqns for } d \text{ features}$$

$$\nabla_b \mathcal{L} = \frac{\partial}{\partial b} \mathcal{L} = \sum S_i y_i = 0$$

$$W = \sum_{i \in C} S_i y_i X_i \quad ; \quad \sum_{i=1}^{n} S_i y_i = \vec{S} \cdot \vec{y} = 0$$

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

$S \in \mathbb{R}^n \in \mathbb{R}$

$n$-values.

"Trick" — $\underline{KKT}$ — Primal Dual form.

Primal Form. } $\min_{\theta, S} \quad \mathcal{L}(X, S, \theta) = \frac{1}{2} \|w\|_2^2 - \sum_{i=1}^{n} S_i y_i (w^T x_i + b)$

$\boxed{k(x_i, x_j)}$

Dual Form } $\max \quad \mathcal{L}_D(X, S, \theta) = \sum S_i - \frac{1}{2} \sum_{j=1}^{n} \sum_{i=1}^{n} S_i S_j y_i y_j (x_i \cdot x_j)$

s.t.

$w = \sum_{i=1}^{n} S_i y_i x_j$ and

$\sum S_i y_i = 0$

$\phi(x_i) \cdot \phi(x_j)$

Pos. Semi-Defn Definition -

° a matrix $\overset{n\times n}{is}$ $\overset{semi}{positive}$ definite if

$v \cdot (A \cdot r) \geq 0$ for any $r^{\neq 0}$ in domain of $A$.

$A$



$v \cdot (Ar) = \|v\| \|Ar\| \cos\theta$

$\theta < \frac{\pi}{2}$

$\theta > \frac{\pi}{2}$

° A kernel fn is pos. semi definite if $0 < 1,3$, and $\overline{K}$ $\overset{pos\ semi\ df}{matrix\ for}$ any input set

⊙ $\mathcal{H}$ = Hilbert space – a complete inner product space

vector space

→ a Hilbert space is a set $\mathcal{H}$

of vectors such that there is a

complete inner product.

dot

limits work properly.

## Spectral decomp. thm:

Suppose $A_{n \times n}$ is pos. def. symm. matrix with eigenvectors $\&$ eigenvalues $\lambda_i, v_i$, $\quad 0 \le \lambda_1 \le \lambda_2 \le \ldots \le \lambda_n$

$$A = \sum_{i=1}^{n} \lambda_i \, v_i v_i^T$$

$v_i \otimes v_i = P_i$

$P_i$ rank-1 projector.

not $v_i^T v_i = v_i \cdot v_i$

- $x_i \cdot x_j$ $\quad\swarrow \; \mathcal{X} = \mathbb{R}^2$ need the dot product between $\underline{x_i} \; \& \; \underline{x_j}$ to do SVM.

- $\underline{K(x_i, x_j)} = \phi(x_i) \cdot \phi(x_j)$

  $\qquad\qquad\qquad \uparrow \;_{\text{corresponds.}} \qquad \int \; \phi: \underline{\mathcal{X}} \to \mathcal{H} \swarrow^{\mathbb{R}^6}$

Gram $\swarrow$-symmetric
matrix-

$$ \underline{K} = \begin{pmatrix} K(x_1, x_2) & K(x_2, x_3) & \text{---} & K(x_1, x_n) \\ \vdots & & & \\ \vdots & & & \\ K(x_n, x_1) & \text{---} & & K(x_n, x_n) \end{pmatrix} $$

KKT

Primal Problem:

$$\begin{cases} \text{minimize: } \mathcal{L}(x,s) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{n} s_i y_i (w^T x_i + b) + \sum_{i=1}^{n} s_i \\ \text{such that: } s_i \geq 0, \forall i \end{cases}$$

Dual Problem:

$$\begin{cases} \text{maximize: } \mathcal{L}_D(x,s) = \sum_{i=1}^{n} s_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} s_i s_j y_i y_j (\vec{x}_i^T \vec{x}_j) \\ \text{using: } w = \sum_{i=1}^{n} s_i y_i x_i, \text{ and } \sum_{i=1}^{n} s_i y_i = 0 \end{cases}$$

The amazing Kernel trick – nonlinear SVM through a kernel and all dot products in the high dimensional space
Done through a kerekl function

$X = \mathbb{R}^2$

Now, we define a kernel $K : X \times X \mapsto \mathbb{R}$, which can take different forms such as:

- Linear kernel: $K(x, \tilde{x}) = x^T \tilde{x}.$
- Polynomial kernel: $K(x, \tilde{x}) = (x^T \tilde{x} - y)^d.$
- Gaussian RBF: $K(x, \tilde{x}) = e^{-\frac{Vert x - \tilde{x} \|^2}{2\sigma^2}}$

$e^{-\frac{\| x - \tilde{x} \|}{2\sigma^2}}$

Consider the polynomial kernel, for $d = 2$, $X = \mathbb{R}^2$, then we have:

$$
\begin{aligned}
K(x, \tilde{x}) &= (x \cdot \tilde{x} + 1)^d \\
&= (x^T \tilde{x} + 1)^d \\
&= (x_1 \tilde{x}_1 + x_2 \tilde{x}_2 + 1)^2 \\
&= x_1^2 \tilde{x}_1^2 + 2 x_1 \tilde{x}_1 + 2 x_2 \tilde{x}_2 + x_1 \tilde{x}_1 x_2 \tilde{x}_2 + 1
\end{aligned}
$$

$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ ; $\tilde{x} = \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix}$

$+ x_2^2 \tilde{x}_2^2$

Kernel is

- ⓪ —
- ① symmetric
- ② pos. semi-def.
- ③ cts. w.r.t. both inputs.

$K(x, \tilde{x}) = K(\tilde{x}, x)$

which interestingly can be re-written in terms of dot product:

$$K(x, \tilde{x}) = (x \cdot \tilde{x} + 1)^d$$
$$= (x^T \tilde{x} + 1)^d$$
$$= (x_1 \tilde{x}_1 + x_2 \tilde{x}_2 + 1)^2$$
$$= x_1^2 \tilde{x}_1^2 + 2x_1 \tilde{x}_1 + 2x_2 \tilde{x}_2 + x_1 \tilde{x}_1 x_2 \tilde{x}_2 + 1 \quad + x_2^2 \tilde{x}_2^2$$

which interestingly can be re-written in terms of dot product:

$$K(x, \tilde{x}) = (1, \sqrt{2} x_1, \sqrt{2} x_2, x_1^2, x_1 x_2, x_2^2) \cdot (1, \sqrt{2} \tilde{x}_1, \sqrt{2} \tilde{x}_2, \tilde{x}_1^2, \tilde{x}_1 \tilde{x}_2, \tilde{x}_2^2)$$

$$\phi : \mathbb{R}^2 \to \mathbb{R}^6$$
$$K : \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}$$

$$K(x, \tilde{x}) = (\phi(x) \cdot \phi(\tilde{x}))$$

$\mathbb{R}^2$

$\mathbb{R}$

$r = x_1^2 + x_2^2$

$\mathcal{X} = \mathbb{R}^2$

$\mathcal{H} = \mathbb{R}^6$

$$\phi : \mathcal{X} \to \mathcal{H}$$
$$x \mapsto \phi(x) = (1, \sqrt{2} x_1, \sqrt{2} x_2, x_1^2, x_1 x_2, x_2^2)$$

$D = \infty$.

$n-1$

$$\phi(x_1, x_2) = (\phi_1(x_1, x_2), \phi_2(x_1, x_2), ..., \phi_6(x_1, x_2))$$

where $\phi : X \mapsto \mathcal{H}$.

$\phi: \mathbb{R}^2 \to \mathbb{R}^6$
$K: \mathbb{R}^2 \times \mathbb{R}^2 \to \mathbb{R}$

Note that $X = \mathbb{R}^2$ is the domain, and $\mathcal{H}$ is the Hilbert space, which is (in machine learning literature) the feature space, and a set of features $\phi_i, \forall i$, is called dictionary.

## Mantra

A major theme in machine learning is that sometimes things actually get easier in higher dimensions !!!.

- A linear plane in high dimensional feature space $\mathcal{H}$, may be a nonlinear curves in the domain space.
- $\mathcal{H}$ is a plane, with calculus with dot products is legit.

The following, we introduce Mercer's theorem, which generalizes spectral decomposition theorem.

$K : X \times X \to \mathbb{R}$

**Theorem 5.5.1 — Mercer's Theorem** generalizer spectral decomp thm.

. Let $K \in L^2(X \times X)$, (i.e. $\int |K(x, \tilde{x})|^2 dx d\tilde{x} < \infty$) such that $T : L_2(X) \mapsto L_2(X)$ by $(T(f)(x)) = \int K(x, \tilde{x}) f(\tilde{x}) d\tilde{x}$ is positive definite. If $\phi_i \in L^2(X)$ is

usual
$Av = U$
eigen
$Aw = \lambda w$

a normalized eigenfunction with eigenvalues $0 < \lambda_1, \leq \lambda_2 \leq \cdots \leq \lambda_N$, Then

exist.

$$K(x, \tilde{x}) = \sum_{i=1}^{N_{\mathcal{H}}} \lambda_i \phi_i(x) \phi_i(\tilde{x}) \tag{5.20}$$

eigen fn.

$T(\phi_i)(x) = \lambda_i \phi_i(x)$
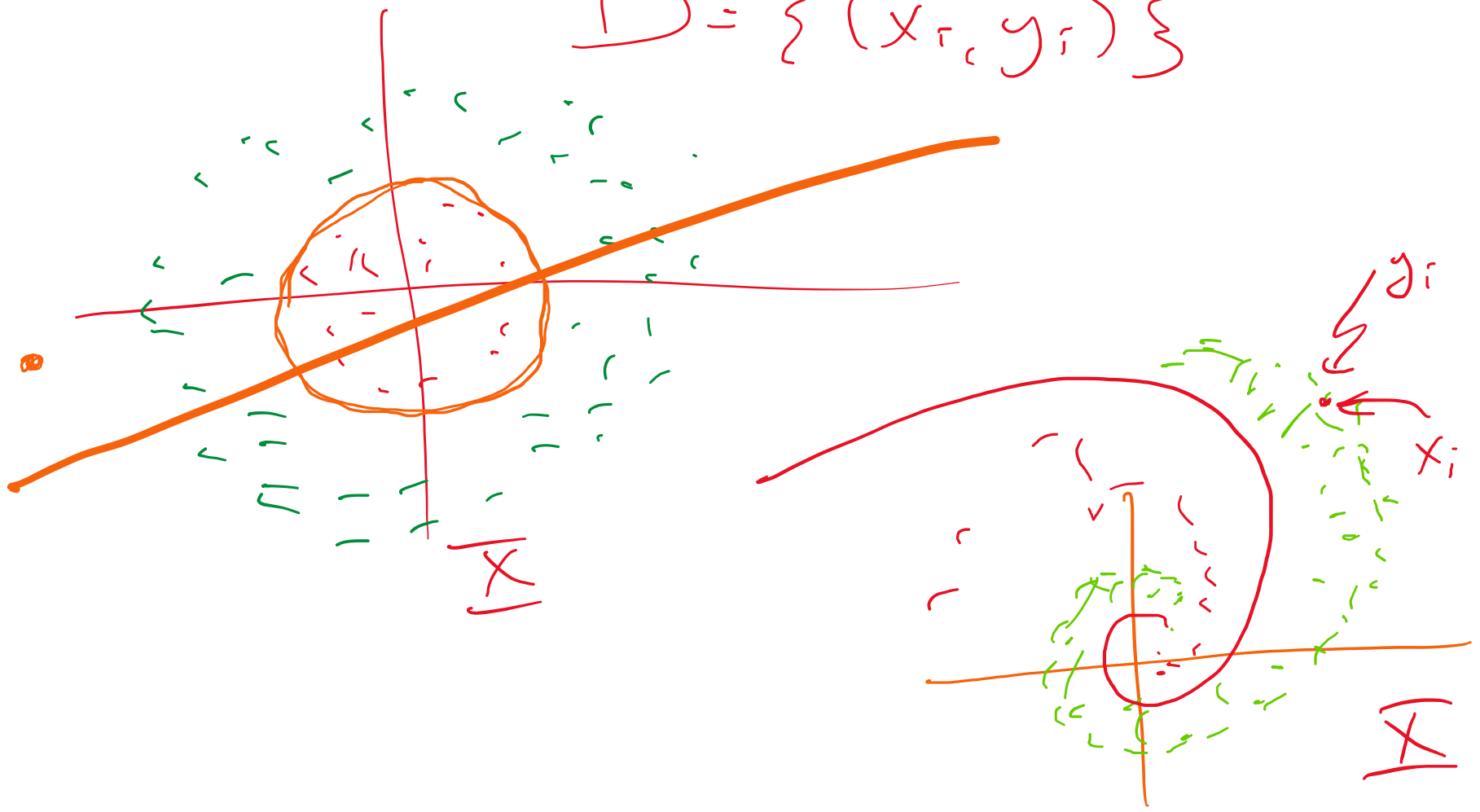
for almost every $(x, \tilde{x})$. Where $N_{\mathcal{H}} = dim(\mathcal{H})$, and the convergance of $K(x, \tilde{x})$ is absolute.

Mercer's theorem itself is a generalization of the result that any symmetric positive-semidefinite matrix is the Gramian matrix of a set of vectors.

$\phi_i$'s exist & I can use them in KSXM.

$$D = \{(x_i, y_i)\}$$



$X$

$y_i$

$x_i$

$X$

And now for something completely different