Pedagogical Lecture: Clustering, SVD, FLD, SVM, concentration of measure

For this lecture we will be looking at images
as points in high (infinite even) dimensional spaces.
It is in fact very natural to think of images as
vectors, whether we are modeling them as functions
or as $n \times n$ arrays of intensities

$$u : [0,1]^2 \to \mathbb{R} \text{ or } \mathbb{R}^3$$

Tasks such as detection, classification, inference
are tasks we would like to do automatically.
And both the Theory and Computation required
are accessible to good undergraduates.

PCA/SVD In this lecture we begin with a classic method for
generating a reduced complexity model of data in high
dimensional spaces: principal component Analysis (PCA)

$$N$$

$$V \equiv n \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \cdots & v_N \\ | & | & & | \end{bmatrix} \leftarrow \text{data, arranged in columns}$$

data $\subset \mathbb{R}^n$

$$C \equiv \frac{1}{(N-1)} V V^T \quad \text{(where we have assumed that the mean has}$$
$$\text{been removed} \cdots \frac{1}{N} V \begin{bmatrix} | \\ | \end{bmatrix} [1 1 \cdots 1] \text{ has}$$
$$\text{been subtracted from V)} \quad \underset{\in N \times 1}{} \quad \underset{1 \times N}{}$$

Now we diagonalize $C$. $\qquad C = Q_c \Sigma Q_c^T$.

Let's use the SVD we met earlier (Lecture 4)

$$\textcircled{1}$$

$$V = \,_n\!\left[ O_L \right]^{\!n} \left[ \begin{matrix} \sigma_1 \sigma_2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{matrix} \;\vdots\; \right]^{\!N} \left[ \begin{matrix} & & \\ & O_R^{\,T} & \\ & & \end{matrix} \right]^{\!N}$$

$$\Rightarrow \quad C = \tfrac{1}{N-1} V V^T = \tfrac{1}{N-1} O_L \, \Sigma \, \Sigma^T O_L^T$$

$$\Rightarrow \quad \text{eigenvalues of } C \text{ are } \sigma_1^2, \sigma_2^2, \dots \sigma_n^2 \text{ from the SVD of } V$$

<span style="color:red">**An exceedingly useful property of the SVD**</span>

$$V_K = O_L \cdot \left[ \Sigma_K \right] \cdot O_R^{\,T} \quad \text{where} \quad \Sigma_K = \left[ \begin{matrix} \sigma_1 \cdots \sigma_K & & 0 \\ & 0 \, 0 & \\ 0 & & \ddots 0 \end{matrix} \;\vdots\; 0 \right]$$

is the best rank K approximation of $V$ in <u>both</u> the $L^2$ operator norm and the Frobenius norm.

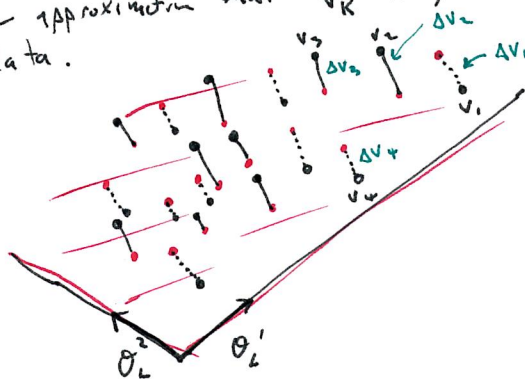<span style="color:red">$\|V\|_F = \sqrt{\sum_{i,j} |V_{ij}|^2}$</span>

<span style="color:red">$\|V\|_2 \equiv \max_{x \neq 0} \dfrac{\|Vx\|}{\|x\|}$</span>

$$\star \quad \|V - V_K\|_F \le \|V - W\|_F \quad \text{any } W \text{ with rank at most } K$$

$$\star \quad \|V - V_K\|_2 \le \|V - W\|_2 \qquad " \qquad\qquad " \qquad "$$

The first inequality says that using the $L^2$ distance, one cannot do better in linear approximation than $V_K$ if you want a $k$-dimensional approximation of the data.



span$(O_L^1, O_L^2)$ gives best 2 dim approx subspace for the data ... the sum of the lengths to the subspace from the data, squared is minimal

$$\sum_i \| \Delta v_i \|_2^2 \text{ is minimal}$$

**Comment:** Because of its ubiquitous usefulness and its power of illumination, students should learn about the SVD as soon as possible.

Thinking about Matrices, Linear Algebra Data with the SVD in your toolset is much easier than the same sans-SVD!

## what ~~**is**~~ is the SVD good for?

① as we have already seen, optimal dimension reductions ... optimal in the $L^2$ sense.

② ... and before even seeing that, we saw it's use in regularization of inverse problems.

③ reduction in computational costs. This is really a corollary or side effect of ①, but it is so important I give it its own bullet. Simply reducing the dimension of the representation can make a huge difference because of the resulting reduction in computational costs.

④ understanding: what are the important components in my data?

## Fisher Linear Discriminant (FLD)

Suppose our data $X_1, X_2, \ldots, X_N$ are each given either the label 0 or the label 1. Assuming that these labeled data points represent some distribution of type 0 n type 1 points well, we can try to build a classifier that separates the labeled data well, believing that it will then be a good classifier for new, unlabeled data.

Data: $\{(V_i, Y_i)\}_{i=1}^{N}$, $V_i \in \mathbb{R}^n$, $Y_i \in \{0, 1\}$

① $\quad m_0 \text{ or } m_1 = \dfrac{\sum\limits_{y_i = 0 \text{ or } 1} V_i}{\#\{y_i = 0 \text{ or } 1\}} \rightarrow N_0, N_1 \quad N_0 + N_1 = N$

② $\quad a \in \mathbb{R}^n$

③ $\quad W_0, W_1 = $ matrices with columns equal to $\{V_i \mid y_i = 0\}$, $\{V_i \mid y_i = 1\}$ respectively.
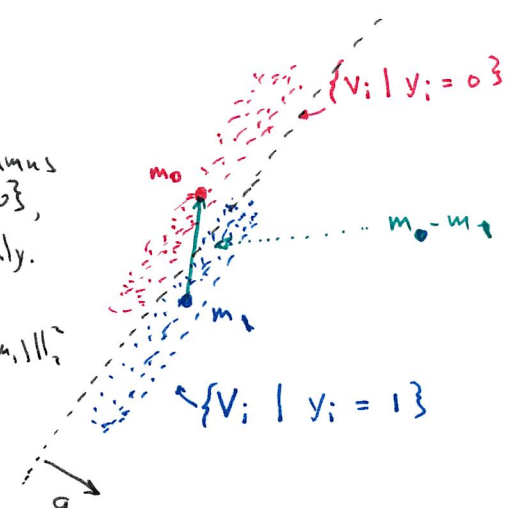
I use

$W_0 - m_0$

to denote

$W_0 - \begin{bmatrix} | \\ m_0 \\ | \end{bmatrix}\begin{bmatrix} 1 & 1 & \cdots & 1 \\ {}_{i1} & {}_{i2} & \cdots & {}_{N_0} \end{bmatrix}$

④ $\quad \sigma_0^2, \sigma_1^2 \equiv \|a^T(W_0 - m_0)\|_2^2, \|a^T(W_1 - m_1)\|_2^2$

and analogously for

$W_1 - m_1$

$\|a^T(W - m)\|_2^2 = a^T \underbrace{(W - m)(W - m)^T}_{S} a$



FLD: "Best" $a$ given by

$$\max_a \; J(a) \equiv \frac{a^T(m_0 - m_1)(m_0 - m_1)^T a}{\sigma_0^2 + \sigma_1^2}$$

$$= \frac{a^T(m_0 - m_1)(m_0 - m_1)^T a}{a^T(S_1 + S_2)a}$$

unsigned

(all we care about is the direction)

Solution: $\quad a = (S_1 + S_2)^{-1}(m_0 - m_1)$

quick proof: $\quad \min \; a^T(S_1 + S_2)a \quad$ subject to $\quad \bar{a}^T(m_0 - m_1) = 1$

$\Rightarrow \quad (S_1 + S_2)a + \lambda(m_0 - m_1) = 0$

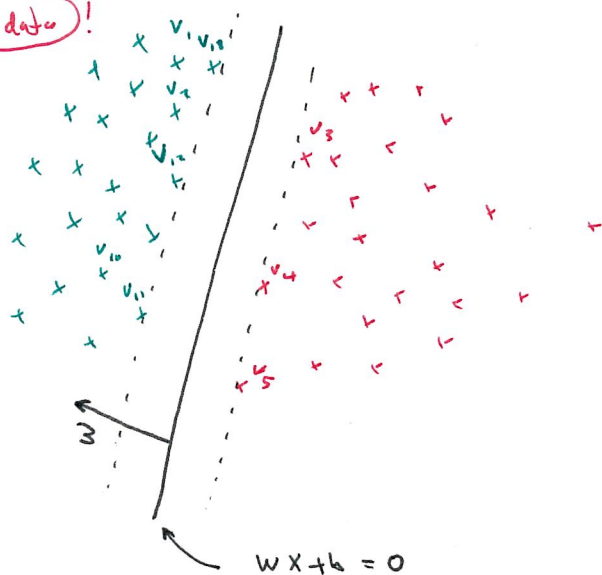$\Rightarrow \quad a = (S_1 + S_2)^{-1}(m_0 - m_1) \quad$ works

Caution: this can work pretty badly!

But: The ideas encountered by students when seeing this are all important... and it is often the first thing you should try with data classification problems.

# Support vector Machines (SVMs)

The next trick in the bag should be support vector machines.

$W X + b = 0$

$$\min \|W\| \qquad \text{subject to} \quad y_i (WX + b) \geq 1$$

$\Rightarrow$ this results in a $W, b \ni$ margin is maximized; margin = min distance from separating hyperplane to data.

Kernel trick: all data is separable!

Idea: map data to a high (infinite) dimensional space where we know how to compute inner products.

why?: The data will, generically, land on the vertices of a simplex $\Rightarrow$ any subset is separable from its complement by a hyperplane in the high dimensional space
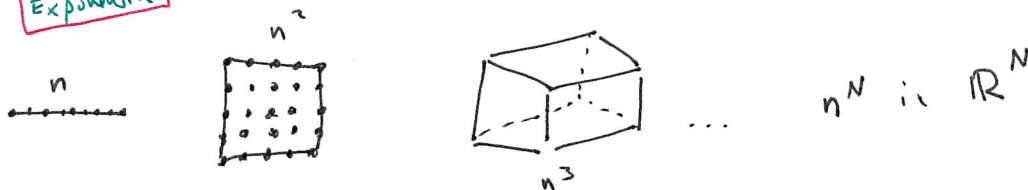


Build the SVM in the high dimensional space.

working high dimensions brings both blessings and curses:

First, the cursing ☺

Exponential



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

... even if you want to "cover" the surface of a cube by only putting points at the vertices, in 100 dimensions you need $2^{100}$ points which is about $10^{30}$ points !!

... and as high dimensional spaces go $\mathbb{R}^{100}$ is very low dimensional !

$C^N$ and $N!$ get huge as $N \to \infty$
fast

but there are blessings !!

Concentration phenomena

These phenomena are both illuminating al entertaining ... students will find this something truly new and surprising ... and accessible to playg around themselves.

We start with volumes of balls: spheres and balls in $\mathbb{R}^n$

$$V(B_r) = \alpha(n) r^n$$

$$V(S_r) = n\alpha(n) r^{n-1}$$

Sterling $\to$ $n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\theta(n)/12n}$

$$n = 1, 2, \ldots$$
$$0 < \theta(n) < 1$$

$$\alpha(n) = V(B_1)$$
$$= \frac{\pi^{n/2}}{\Gamma(\frac{n}{2}+1)}$$
$$= \frac{\pi^{n/2}}{\frac{n}{2}\Gamma(\frac{n}{2})}$$

(Playing around with spheres and balls in $\mathbb{R}^n$ can teach a great deal actually)

(1)

notice:
$$V(S_r) = \frac{\partial}{\partial r} V(B_r)$$

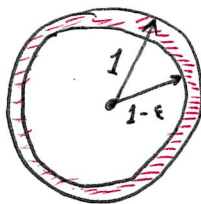what is $\frac{\partial}{\partial r} V(S_r)$?

answer:
$$\frac{\partial}{\partial r} V(S_r) = n \cdot n-1 \, \alpha(n) r^{n-2}$$
$$= \frac{n-1}{r} \left( n \, \alpha(n) r^{n-1} \right)$$
$$= \frac{n-1}{r} V(S_r)$$
$$= \int_{\partial B_r = S_r} \vec{H} \cdot \vec{n} \, d\sigma = \int_{S_r} \vec{H} \cdot \vec{n} \, d\mathcal{H}^{n-1}$$

total mean curvature of sphere of radius r

(2)



$$\frac{V(B_{1-\epsilon})}{V(B_1)} \xrightarrow[n \to \infty]{} 0 \qquad ! \qquad \text{(for any } \epsilon > 0 \text{ no matter how small)}$$

i.e. $(1-\epsilon)^n \xrightarrow[n \to \infty]{} 0$
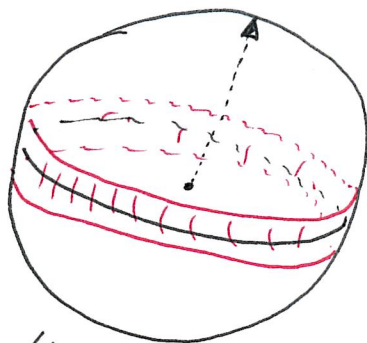
As $n \to \infty$ almost all the volume of a ball is in a very thin shell on the boundary!

This is the essence of concentration phenomena

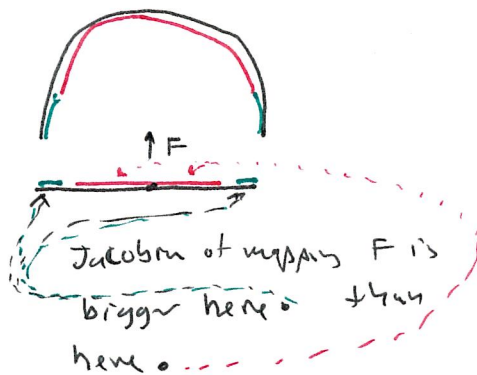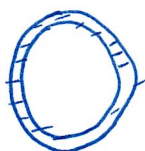③ <u>Equitorial Bands on Spheres</u>



contains almost all the area on the surface as $n \to \infty$ for any $\varepsilon > 0$ !

<u>Explains</u> : random vectors in high dimension are almost orthogonal.

get hemisphere from warped disk



Jacobian of mapping F is bigger here. than here.

④ Tons of other things to explore

* Gaussian dist in $\mathbb{R}^n \Rightarrow$ most vectors have length $\sqrt{n}$ when $C = [\dot{\,}, 0]$

* volume of unit Ball $B_1 \to 0$ $n \to \infty$

* cubes and balls intersect ...

⋮

* Lipschitz functions on spheres are approximately constant as $n \to \infty$ (i.e. on a set of large measure function it is within $\varepsilon$ of a median value).

∴