

Kernel Methods :-

Kernel - PCA :-

1 → Given a set of data x^i be "invent" a function ("Kernel")

$$K(x^i, x^j)$$

→ The only requirement on this funcⁿ is that it should be symmetric

→ Sometimes, we need it to be two-defined

$$\text{i.e., } K(x^i, x^j) = K(x^j, x^i)$$

→ Then we compute the $N \times N$ Kernel matrix

2 → $K^{ij} \rightarrow g^{ij} = K(x^i, x^j) + \gamma_{ij}$ (we perform double-centering to take our kernel matrix to a gram matrix)

$$g^{ij} = -\frac{1}{2} (k^{ij} - \frac{1}{N} \sum_k (k^{ik} + k^{jk}) + \frac{1}{N^2} \sum_{lm} k^{lm})$$

double centering

→ We can prove easily that whatever k we take, assuming its symmetric, the resultant (g^{ij}) must be a Gram matrix.

scalar product so not two defined

i.e. $\begin{aligned} g^{ij} &= g^{ji} \quad (\text{symmetric}) \\ &\leq g^{ii} = 0 \quad \forall i \end{aligned}$

i.e., resultant g^{ij} satisfies the properties of a Gram Matrix due to double-centering.

3 → We perform MDS with $\underline{g^{ij}}$

will simply give us a list of each vector.

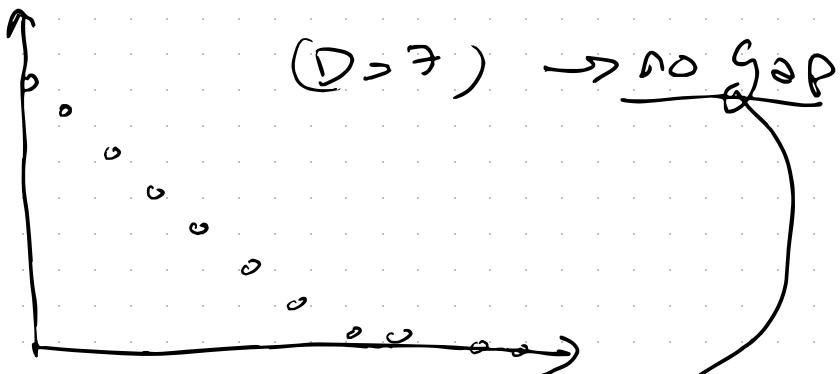
→ Now how do we get K^{ij} ?

① Let's say a null hypothesis / ground state x^j , $K^{ij} = x^i, x^j$

→ If we do this, we are basically doing PCA.

→ Of course it doesn't make sense to do this; easier to diagonalize C than S .

→ But spectrum would be like (for eg.):



So, we are not happy with our choice of kernel.

① So, let's choose a new kernel

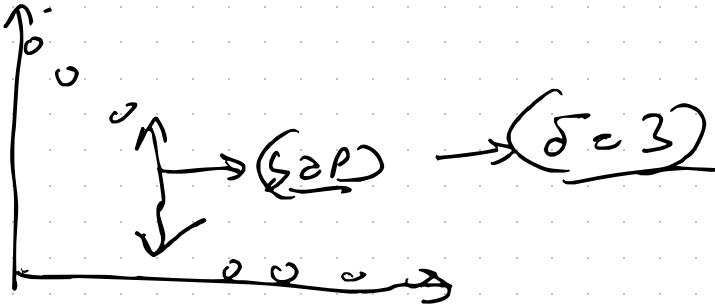
$$K^{ij} = x^i \cdot x^j + \alpha (x^i \cdot x^j)^2$$

We compute spectrum again

α = variational parameter

S. Matk, $\alpha = 0 \rightarrow \text{PCA}$
 $\alpha = 1$

(for eg)



∴ we have found a kernel when 3 eval
of \mathcal{L} explain 100% of the variance.

① TL;DR: We choose a K-junc²
in such a way that the
spectrum can be truncated as
early (for as few eigenvalues) as
possible, i.e., spectrum of gram matrix
can explain all / target variance for
as few eigenvalues as possible.

⇒ For eg, $K_{a,b}(x^i, x^j) \equiv$ kernel func²
with variation parameters
(a) + (b).

$(x_{a,b})$ spectrum
must be obtained

Find (a, b) s.t.

+ target
embedding
dimension, can be
as large as possible
from ③

$$\sum_{\alpha=1}^3 \lambda_{\alpha}(a, b)$$

$$\sum_{\alpha=1}^D \lambda_{\alpha}(a, b)$$

i) as large as possible

→ This means we have found a specific kernel that allows a significant dimensional reduction.

maximize

$$\sum_{d=1}^D \lambda_d (a, b)$$

$$\frac{\sum_{d=1}^D \lambda_d (a, b)}{D}$$

(δ is supplied)

→ We introduce the kernel s , that we can estimate the distances in a manner which is compliant with the features of the embedding manifold.

→ How do we choose the kernel?

- ① Kernel is rotated to a suitable dist. on the manifold.
- ② (Traditional idea) The kernel defines implicitly a scalar product in an extended feature space.

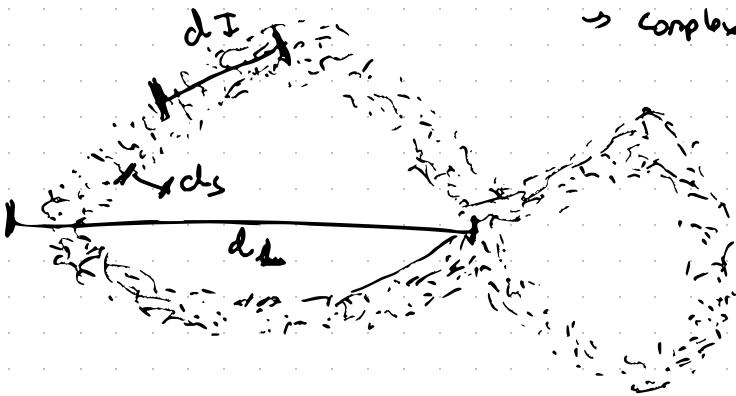
→ Our data manifold is nonlinear (curved) so when we add nonlinear features to represent our data.

(Nonlinear features
= = extended feature
space)

(we extend feature
space from ordinary
Cartesian coords to
add $\sqrt{x_1}, x_1^2$ etc.
& other nonlinear features)



Traditional idea - \rightarrow noise on small scale
 \rightarrow curvature on medium scale
 \rightarrow complexity on topological scale



$d_s \equiv$ Small length scale

$d_L \equiv$ large length scale on which there is
topological complexity,
curvature etc.

(intermediate)

$d_I \equiv$ Sweet spot length scale

\equiv locally flat manifold is on a
hyperplane in this region

so only on this length scale it makes sense to
approximate the manifold with a hyperplane.

→ i.e., Simple way of defining a kernel :-

LOSEL irrelevant $|x^i - x^j| < d_S$

∴ we assume fluctuations at this scale are noise

FAR, irrelevant $|x^i - x^j| > d_L$

∴ if 2 data pts. are so far away, we can't model them so we disregard them,

ACCEPTABLE :-

$$d^{ij} = 1 - \exp\left(-\frac{\|x^i - x^j\|^2}{2 d_I^2}\right)$$

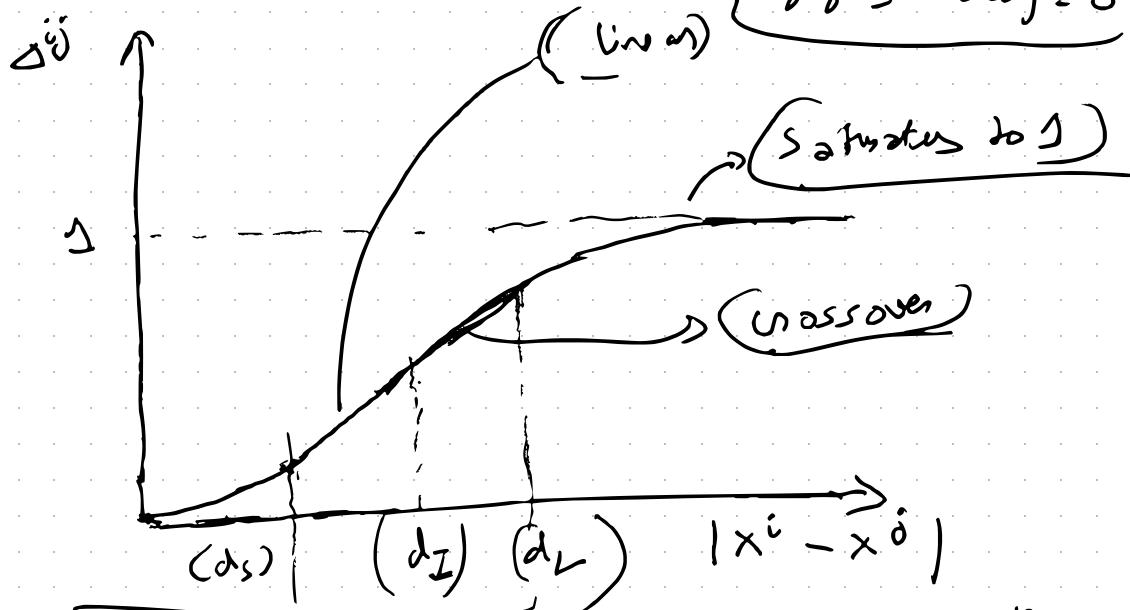
(define \geq possible length scale)

(Gaussian Kernel) $\left\{ \frac{"1"}{\cdot} \text{ not relevant close to double centering} \right.$

k is
 \approx Gaussian

(we can keep it, or remove it)

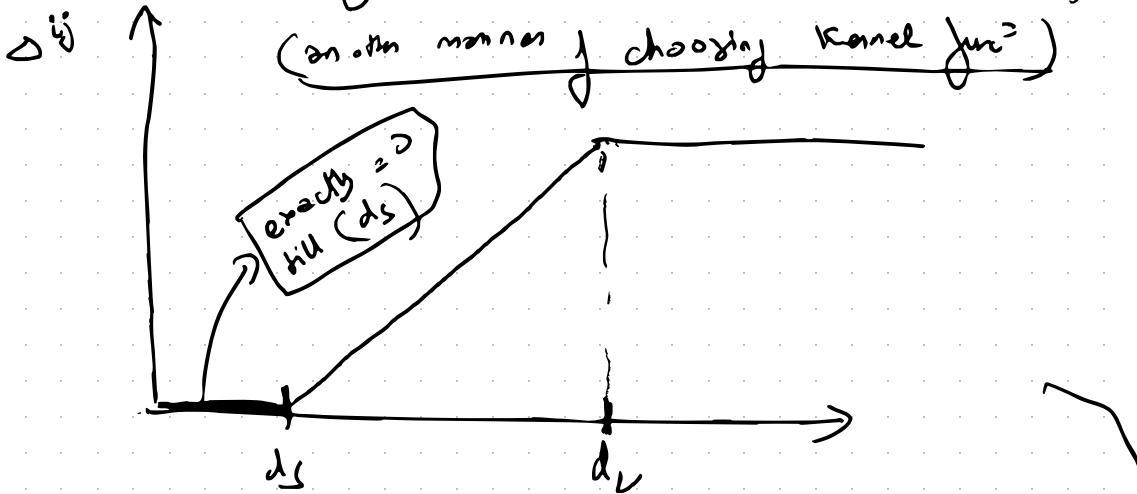
→ all pts. with $\text{dist} < d_s$, are lumped into something negligible, basically ≈ 0



→ pts. that are very far, are all lumped into the same dist., implying "far", here chosen = 1

→ When we use a Gaussian kernel, we only have one parameter, the variance, so we sort of implicitly define distances on our data manifold, in which the only relevant distances are of order (d_I). So, we have just a single parameter, (variance).

• SKETCHMAP :- Case of Kernel-PCA :-
 (an often manner of choosing kernel func²)



- another possible choice of kernel func²
- This is called "SKETCHMAP"

Principle :- if we have a data manifold that is affected by noise on short - scale & on the large scale is affected by complexity, we only have a suitable intermediate dist. range (d_I) on which we are interested, then the sketchmap is

→ this suitable kernel automatically —
 selects the intermediate dist. for our feature learning / dim. reduction.

◻ Let's take a step back to MDS :-

→ We find \mathbf{Y}^{u} , s.t., $\|\mathbf{Y}^{\text{u}} - \mathbf{Y}^{\text{d}}\|$ $\approx \mathbf{S}^{\text{ij}}$

→ Using kernel methods for features that we "learn" will reproduce ONLY \mathbf{Y} instead of dist.

→ If we use MDS algo., we solve this problem analytically.

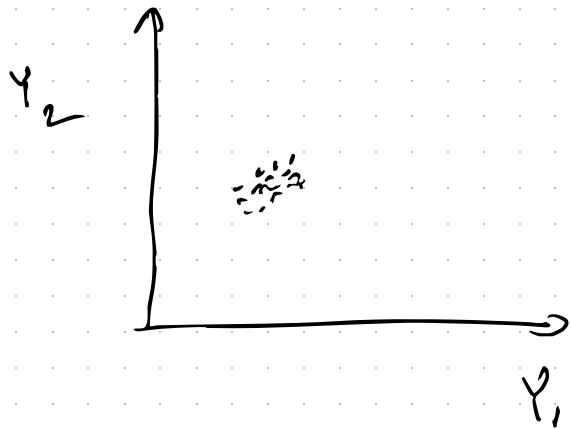
→ we have a lot of variational freedom,

→ In a very complex topology, if we say everything farther than a dist. is the same, then we lose information for sure; we can't recover the original topology.

→ getting a dim.

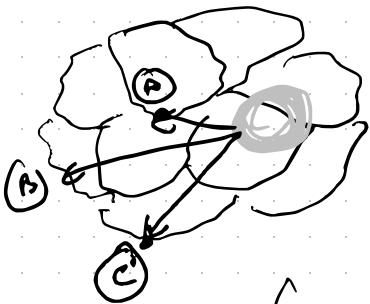
new dim. that makes sense is much more likely.

→ This is because $(\mathbf{Y}^{\text{u}})^{\text{s}}$ don't even attempt to reproduce properties of our system on the large scale.



- is a reduced dim. representation,
- points that are close in $(Y_2 - Y_1)$ space
are not necessarily close in original dataset. \therefore NN's in Y are not necessarily NN's in X & vice-versa.
- Because if they are closer than a threshold,
all the dist. are the same.
- i.e. we have a lot of variation in freedom
to obtain the interm. distance.

→ Imagine a complex topology on a large scale.



→ imagine we do MDS on the many of the grey doughnuts.

→ anything smaller than it is considered noise, everything outside is too far & also ignored.

→ In our low dim. projection, our results will not look like our original manifold.

↓
C :: dist. to A, B, C is treated as exactly the same)

↳ ; at a certain scale, it is locally compliant with our dataset, but at the large scale, it's not.

→ ; we can measure dist. in $\{y^i\}$ space, but these dist. are relevant only in (rid) range.

TL; DR:-

(any funcⁿ that is symmetric in its arguments)

→ $k(x^i, x^j) = k(x^j, x^i)$ is our only requirement.

$$k(x^i, x^j) = 1 - \exp \left(-\frac{\|x^i - x^j\|^2}{2\sigma^2} \right)$$

(Gaussian Kernel)

Concept of suitable length scales

(doubt from before)

26/11/23

Q. We derived double-centering formula considering
 Δ^{ij} as euclidean dist. Then why does it
work with kPCA & ISOMAP?

A.

$$\Delta^{ij} = \|x^i - x^j\|^2$$

$$g^{ij} = -\frac{1}{2} \left(\Delta^{ij} - \frac{1}{n} \dots \right)$$

(double centering formula)

Then $g^{ij} = x^i, x^j$
(exactly)

generic
form
of the
coors.

In general, if $\Delta^{ij} \geq \Delta^{ji} \Rightarrow \Delta^{ji} = \Delta^{ij}(x^i, x^j)$

(i.e., it only satisfies
symmetry operation
w.r.t. the indices, (i, j)]

(most general
case)

then $\rightarrow g^{ij} \neq x^i, x^j$
(in general)

however $\sum_j g^{ij}$ with satisfy the properties

$\sum_j g^{ij}$ Gram Matrix is i.e. :-

$$\sum_j g^{ij} = 0, \forall i$$

→ This property allows us to interpret

g^{ij} as a scalar product between (z^i, z^j)

which are other, newer features,

not (x^i, x^j)

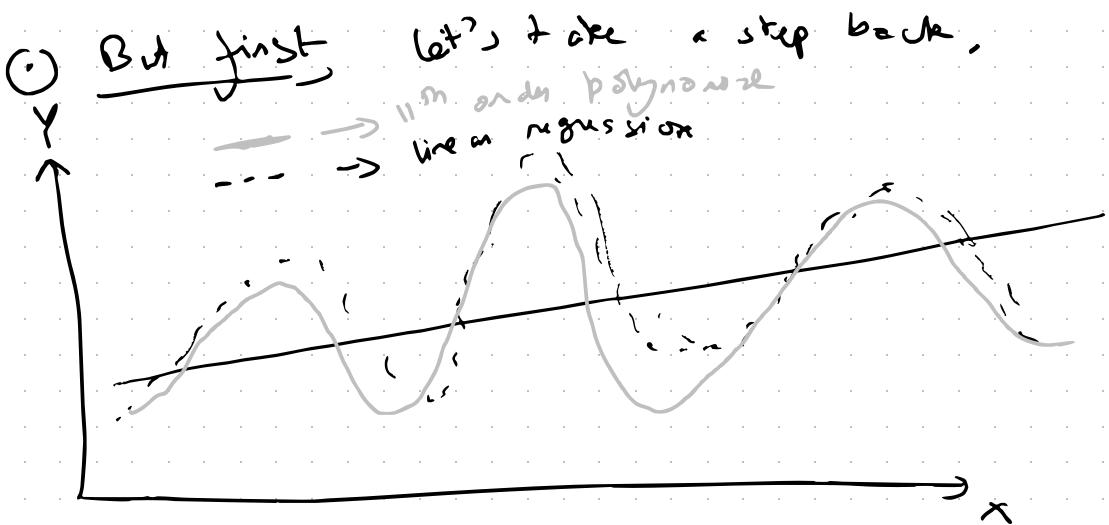
$$\therefore g^{ij} = z^i, z^j \quad z^i \neq x^i \\ z^j \neq x^j$$

→ P.B. ; g^{ij} can be assumed to be derived from a kernel.

\rightarrow i.e. we can interpret the Gram Matrix as a scalar product between centered features in an extended feature space.

(simply because it's symmetric)
 $\forall \sum_j g^{ij} = 0$, if i is satisfied)

This leads us well into the topic for kernel tricks



$X = \text{master variable}$

$Y = \text{depends on } X$

→ imagine we try to do linear regression

↪ i.e., we assume $y^i = a + b x^i$

if we try to determine $(a) + (b)$

s.t. that a straight line fits the data.

→ But, this gives us a very bad result & is not a good model, ∴ using linear regression with the original features (just x^i) → we get really bad results.

→ So how do we improve?

let's say instead of a linear funcⁿ we try to represent it as a polynomial to increase the basis set.

$$y^i = a + b x^i + c x^{i^2} + \dots$$

→ so if we extend this "basis" enough, eventually we can reproduce the funcⁿ. (let's say if we include 12 terms)

↪ this is still a linear fit.

- \Rightarrow even for a very complicated func²
 we can reduce it to a linear regression
 problem
- \Rightarrow That is $\overrightarrow{\text{it doesn't matter how complex the}}$
 $\text{func}^2 \text{ is if we except the cost of extending}$
 $\text{in feature space introducing nonlinear func}^2$
 $(x^2, x^3, \sin(x), \cos(x) \text{ etc.})$ then we can
 basically reproduce any possible func² by linear
 regression.
- \Rightarrow In other words $\overrightarrow{\text{any nonlinear twisted data}}$
 manifold can in principle be represented by
 performing a suitable representation in a
 feature space which includes a sufficient number
 of nonlinear func² of the original coordinates $\overrightarrow{\text{}}$
- \Rightarrow We can also build more complex features
In this case $\tilde{x}^i = y^i$
 & now we can build features which
 are of (x^i, y^i) & still do this
 process with even more complex feature space.

→ In practice, if we accept to extend our feature space into N dimensions (with non-linearity), then we can map any data representation into a linear problem.

Statement :- If we extend our feature space introducing extra features which are non-linear combinations of our original features (x^i) then it is more likely that we will be able to represent our data manifold with a linear regression.

- For eg, if the manifold is a hyperplane then original features are far enough sufficient to describe it with linear regression, which is just PCA
- If the manifold is more complex & NOT a hyperplane, then adding non-linear features makes it more likely that we can "linearize" our problem.

\Rightarrow So what exactly, in mathematical terms,
does this mean?

- ① \rightarrow Let's say we are in 2-D.
 \rightarrow original feature space: (x_1^i, x_2^i) for each data point (i)

→ extended feature space :-

$$(x_1^2, x_2^2, x_1 x_2)$$

so we pass from = representation
 on data in $D=2$ to 2
 $\sim \sim \sim D=5$

→ we can do PCA in this extended feature space. This way we can find if there are any new emergent linearities and as a result of extending the feature space,

→ A (g^{ij}) derived from a kernel by
Gram matrix double centering can be interpreted as a scalar product in an extended feature space.

→ Virtually any Gram Matrix built by double centering from a symmetric kernel can be interpreted as a scalar product in an extended feature space.

⇒ Eg :-

① Polynomial Kernel :-

$$k(x^i, x^j) = \sum_{\alpha} p_{\alpha} (x^i \cdot x^j)^{\alpha}$$

(coefficients)
(dot product - scalar product)

$$\rightarrow \text{if } P_1 = 1, P_{d+1} = 0$$

$$K^{ij} = x^i \cdot x^j$$

(i.e., kernel is just the standard scalar prod. b/w the original features)

$$\rightarrow \text{if } P_1 = 1, P_2 = 1, P_{\text{rest}} = 0,$$

$$K^{ij} = x^i \cdot x^j + (x^i \cdot x^j)^2$$

'if we have feature space = L :-'

$$\text{then } (x^i \cdot x^j)^2$$

$$= (x_1^i x_1^j + x_2^i x_2^j)^2$$

$$(x_1, x_2) = (x_1^i)^2 (x_1^j)^2 + (x_2^i)^2 (x_2^j)^2$$

$$+ 2 x_1^i x_2^i x_1^j x_2^j$$

$$(x_1^2, x_2^2), \sqrt{2} x_1 x_2$$

∴ our new feature space :-

$$x_1, x_2, z_3 (= x_1^2), z_4 (= x_2^2), z_5 = \sqrt{2} x_1 x_2$$

$$\therefore (x^i, x^j)^2$$

$$\Rightarrow z_3^i z_3^j + z_4^i z_4^j + z_5^i z_5^j$$

$$\begin{aligned}\therefore k^{ij} &= x^i, x^j + (x^i, x^j)^2 \\ &= z^i, z^j\end{aligned}$$

to show

$$z = (x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2)$$



→ extended feature space

→ z is the new basis &
 $(x_1, \dots, \sqrt{2}x_1 x_2)$ are components of the vector

→ even for higher powers, $(x^i, x^j)^3$

etc. we can always group together
the terms that depend on (i)
& those that depend on (j) .

⇒ ∵ $K^{ij} = z^i, z^j$ is always
preserved as long as (z) is
defined in a specific manner,
only that it will probably contain much
more complex features.

⇒ ∵ $K^{ij} = \sum_{\alpha} P_{\alpha} (x^i, x^j)^{\alpha}$

is a scalar product in an
extended feature space.

→ we can extend $\alpha \rightarrow \infty$
basically & can choose (P_{α})
in a specific manner.

\rightarrow if we do this

$$k^{ij} = f(x^i, x^j)$$



can be any generic func² that can be Taylor expanded. It can be an exponential/Gaussian/Lorentzian etc.

∴ This can be thought of as a scalar product in an extended feature space.



This is called the "Kernel Trick"

Kernel Trick :- defining a generic, symmetric func² of my original features & Taylor expand this func² & this expansion is implicitly the def² of an extended feature space. i.e. if the Kernel is an exp(), then the feature space can be considered infinitely large.

- ⇒ This is why kernel methods are extremely powerful for regression, dim. reduc², & feature learning.
- ⇒ if we represent distances with a specific nonlinear kernel, we extend the data representation in a very high dim space. Then in principle, linear regression becomes meaningful.

For eg,

$$K(x^i, x^j) = \exp\left(-\frac{\|x^i - x^j\|^2}{2\sigma^2}\right)$$

$$= 1 - \frac{\|x^i - x^j\|^2}{2\sigma^2} + \frac{1}{2\sigma^2}$$

↓
after Taylor exp, we can't write it as a scalar product without double centering

→ Some terms can be written as a

scalar pdt.

$$\|x^i - x^0\|^2$$

$$2x^i$$

& some
cannot

be written

as a scalar

pdt.

$$x^i x^j - 2x^i, x^j + \|x^j\|^2$$

not a

scalar pdt.

written as
a scalar pdt.

⇒ TL; DR : a generic Gaussian Kernel will
NOT be centered.

∴ The features that we try to introduce
must be centered.

→ Summary :-

(i) perform PCA in an extended feature space (can even be ∞ -sized feature space)

(ii) diagonalize the Gram Matrix;
Covariance matrix in ∞ -feature space
can't be diagonalized.

(iii) This allows us to perform linear reg./PCA
in this extended feature space.

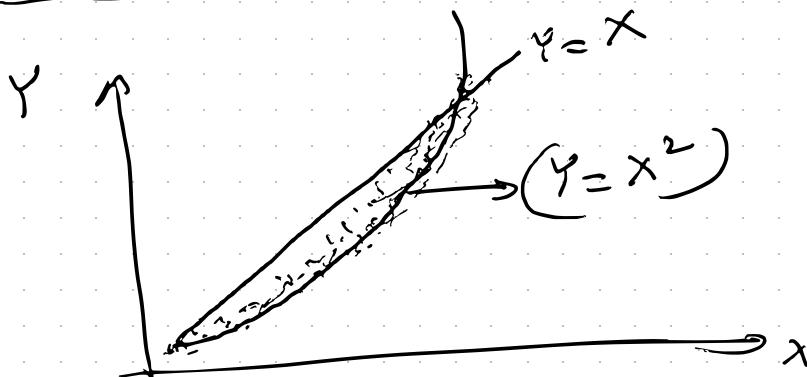


"Kernel Trick"

(iv) Problems :- This kernel trick may not
always be useful.

In order to be really useful, we have
to select our kernel in a really
smart manner.

→ Let's go back to the problem of regression :-



→ Having feature space only with (X) is
a bad result.

→ extending feature space with (X^2) does
very well. (\because we have a quadratic funcⁿ)

$$\therefore \underset{J}{\underbrace{Y}} = a_1 X + b_1 X^2$$

$$a = 0$$

$$b = 0.3 \text{ (for eg.)}$$

→ so if both the features have same weight from our manifold is NOT linear.

→ so basically here our manifold is only linearly be passed away (x) & keep only (x^2).

→ if we can do linear regression using a basis where we have ($1, x^2, x^3, \dots$) but the coefficients of this linear regression are non-trivial $\rightarrow (0, 0.7, 0, 0, \dots)$

⇒ if we choose a generic kernel this is like representing our data on a basis where the coefficients of the data are fixed.

→ in the polynomial kernel of order (2)
we represent our data in the extended
feature space: $x_1, x_2, x_1^2, x_2^2, \sqrt{2}x_1 x_2$

→ if we use this kernel & do PCA, we
will not get rank = 1 matrix

(on the
quadratic
data of the
prev. page)

In order to get
rank = 1 matrix we
must remove all features
except for (x_1^2) .

↓
∴ even in a completely trivial case, while
we can do PCA in extended feature space
by introducing nonlinear function. However this
doesn't mean we transform our data into
a hyperplane. So to get a hyperplane we must
get rid of useless features.

$\Rightarrow \therefore$ There is no one generic good kernel. So given our data manifold we should tried & error many diff. kernels & settle with the one that gives us the eigenvalue spectrum that can be truncated earliest.

II

(Subtle trick)

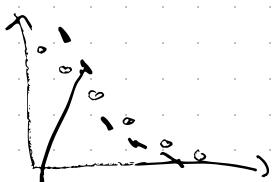
(in general)

→ this is for manifold learning, but also applies to kernel regression (kPCA).

→ TL;DR : ① choose a kernel

② compute the spectrum

③ settle on kernel with best performance (captures most variance, earliest +)



(-) better kernel

$$\rightarrow \sum_{\alpha=1}^n \lambda_\alpha = D$$

(n chooses dimensions out of D)

$$n < D$$

Should be as large as possible

The kernel for which this is max. is the better kernel.

For a Gaussian Kernel

$$K() = \exp \left(-\frac{\|x^i - x^0\|^2}{2\sigma^2} \right)$$

single unidirectional parameter.

What is correct (σ)?

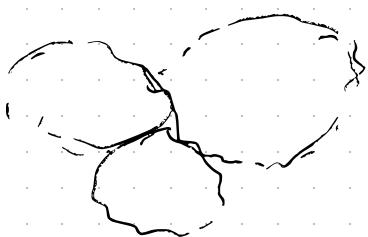
That which makes (ϕ) as large as possible for a fixed (n) sum of components

→ For polynomial kernel we have many more variational parameters to play around with.

N.B., The own. value for (n) is $= ID$.

But it is much safer to set (n)

to larger than (ID) . We will explore this more in the lectures for ID .



→ If this is our data manifold even if $(ID = 1)$ it is completely impossible to represent it globally in $\underline{1-D}$. So,

n must be $> ID$ to be safe. i.e. Computing ID is

a good guideline to choose (n) but it is just a lower bound.

"sufficiently complex topology makes it
safe to choose $(n > ID)$



Graph Layout Projection Methods:

- These methods come from graph theory.
- Aim: find best low dimensional representation of a graph.
- In ISOMAP we build a neighborhood graph & a possible manner of proceeding is computing geodesic dist & then performing MDS.
So, if we have our features we build a neighborhood graph & then apply ISOMAP.
- ⇒ However there are some methods that were taken from graph theory & have become very popular in explicit dimensional reduction.

→ There are 3 different methods:

(1) SNE (Stochastic Neighbourhood Embedding) → Hinton, Roweis
(2003)

(2) t-SNE (t-distributed SNE)
; Hinton et al

(3) UMAP :- (McInnes, Healy, Melville)
(2018)

new standard for producing
2D graphs from very high dim.
data

$\rightarrow x^i \equiv$ original feature space

\rightarrow instead of building a neighborhood

graph in a deterministic manner like
in Isomap, $\xrightarrow{\text{two nodes}}$

\rightarrow we assume links of the graphs are random
variables which exist or not with a
probability that depends on their
neighbors

$\rightarrow \therefore \exists \text{ link between } (i) \text{ AND } (j)$

\hookrightarrow either 0 (doesn't exist) or 1 (exists)

$\xrightarrow{\text{binary event}}$

$$\left(\sim \frac{k_{ij}}{\sum_j k_i} = p_{ij} \right) \xrightarrow{\text{characterized from a binomial distribution}}$$

→ more explicitly

→ prob. to observe link b/w (i) & (j)

given $(j) = P_{i|j}$

→ $\sum_j K^{ij} \equiv$ normalizing over the
 2^{nd} index

$$\therefore P_{i|j} = \frac{K^{ij}}{\sum_j K^{ij}}$$

\therefore these events are characterized
by the kernel (K^{ij}) that

is normalized & harvested from
a binomial prob. distrib.

→ simple choice for the kernel is Gaussian

$$\therefore K^{ij} = \exp\left(-\frac{\|x^i - x^j\|^2}{2\sigma^2}\right)$$

→ One of the biggest differences b/w σ^2

SNE / t-SNE / UMAP is how

σ^2 is chosen.

→ In UMAP

σ is dependent on the data points & is somehow dependent on the distance from the (k^{th}) nearest neighbours of (i) .

$$\therefore \sigma \rightarrow \sigma_i \approx \|x^i - x^{k \text{ N.N. of } i}\|$$

→ i.e. end up with a Gaussian kernel whose width depends on how dense is our neighbourhood.

→ i.e. we get P_{ij} (prob. of existence of links) from this kernel, it means the data points that are closer than k N.N. have a high prob of = link.

→ in datasets that are farter from the
~~farther~~ NN, the prob of a link decreases
 exponentially.

→ In all these 3 methods the general
 idea is to define another kernel in
 a feature space of smaller dimensions.

$$R^{ij} = \exp\left(\frac{-||y^i - y^j||^2}{2\sigma^2}\right)$$

(Gaussian)

\downarrow

(t-SNE)

→
 lives in a lower
 dim. space defined
 by coords. $\{y^i\}$

$\frac{1}{1 + ||y^i - y^j||^2}$

(t-SNE)

$\xrightarrow{\text{t-SNE}}$ t-Student
 distribution

Cum AP ← { $\frac{1}{1 + a \cdot ||y^i - y^j||^{2b}}$ }

→ uMAP has 2 extra metz parameters ($\alpha + \beta$) that determines how our activity decays.

→ In all these 3 cases

(i) $x = p_{ij}$ (joint prob. from condⁿ prob.)
(prob. in original high dim. space)

(ii) $y = \pi_{ij}$ (built with m F kernel)
(prob. in lower dim. space)

→ we want $p_{ij} \rightarrow \pi_{ij}$:-
→ i.e. we minimize the (KL) divergence
to have these 2 sets of prob. as close to each other as possible to justify our dim. reduction.

$$KL(P \parallel \pi)$$

$$= \sum_{ij} P_{ij} \ln \frac{P_{ij}}{\pi_{ij}}$$

(KL divergence) [explained more
in NN class]

→ Let's find $\tilde{\pi}$ such that

$$\pi_{ij}(P^{\tilde{\pi}}) = P_{ij}, \ln 1 = 0$$

$\therefore KL(P \parallel \tilde{\pi}) = 0$

→ $KL(P \parallel \pi)$ is positive-defined &
lower bound is $= 0$, \therefore to minimize
 $KL(P \parallel \pi)$, we find the probs $(\pi^{(i)})$ that
match as close as possible with $(P^{(i)})$.

$\rightarrow \text{TL; DR}_2$

i) your data is connected by stochastic links.

ii) these links are thought of as harvested from a prob. distrib² (chosen by us, typically Gaussian)

iii) then we try to reproduce the original prob. with a model that is functionally identically but in a lower dim. space.

iv) we have freedom to choose the kernel & diff. methods differ in the exact functional form of these prob.



→ So far we have seen ↗

② PCA

① MDS

↳ ISOMAP
kernel PCA

③ t-SNE (Stochastic Neighborhood Embedding)

↳ t-SNE
UMAP

⇒ Desired qualities of UDL techniques:

(1) significant dimensional reduction.

$$\mathbb{R}^D \rightarrow \mathbb{R}^S, (S < D)$$

$S \geq ID$ (intrinsic dimension)

→ If we want/need data visualization,
we need $S = 2 / 3$, but really
even $S=3$ is a problem, $S=2$ is
ideal.

→ If $ID > 2$, it'll be impossible
to obtain a good visualization of our
data.

(2)

In the case that we want to use our reduced dim. representation for further analysis \rightarrow typical case



The neighbourhood of reduced representation should NOT be inconsistent with original neighbourhood



That is points close in X, should also be close in Y,

\rightarrow This requirement can be at odds with topology
e.g., if all pts. are on a sphere & we want to represent it on a plane.

\rightarrow Such an attempted representation would violate this requirement.

(3) Interpretability of Y coordinates.

In PCA, Y vars are linear combos of original variables, so we can interpret that our PCs are dominantly influenced by 50 & 50 features.

In other methods, Y coords are either very complex funcs of original coords, or like in SNE they are stochastic variables learned on the original data. So in SNE interpretability is garbage.

	0	1	2	3
PCA	x	x ↪ ✓		✓
MDS - ISOMAP	✓	— ↪ —	x ↓	
MDS - RP PCA	✓	x ↪ ✓	x ↓	
SNE - tSNE	✓	✓	x	x

$\times \downarrow$ = doesn't have the property but to a lesser extent than (x)

For eg, RP PCA for (3) is hard to interpret unless specific results that allows us to understand something. NOT as hard as SNE.

→ There is also a $\overset{\circ}{\partial} M$ desired property.

→ (1) It should work good on complex, non-trivial manifolds (which are very curved & twisted).

→ (1) & (3) are sort of inversely correlated

→ In PCA & kPCA :-

→ (1) & (2) are anticorrelated.

→ (2) maintains topology & if we want a very strong dim. reduction, we destroy the topological properties.

→ So, if (2) is good [\because PCA & kPCA has great topology preservation], (1) is bad & vice-versa.

→ If we include # of variables that describe our variance sufficiently, (2) is preserved, we only do a rotation of our representation & throw away something which is small, so (2) is well maintained.

→ ISOMAP → suitable for preserving topology, \therefore we unfold the manifold.

We observe the same anticorrelation behaviour as in PCA & kPCA.

- if we perform ISOMAP on a spherical manifold, we are in serious trouble if we unwrap sphere into something \cong 1D
- ⇒ if we use a representation by ISOMAP \cong 1D we can still find a reasonable set \cong (hence \hookrightarrow symbol in the table)

SNE \rightarrow t-SNE, By def² (1) is great,

" we usually use it for visualization (by def²), of course this ensures (2) is horrible. These methods are states of the art for obtaining 2D representations

⇒ Except for in PCA, (3) (interpretability) is horrible for all other methods. This is because (3) is not (explicit & differentiable) function of origin of x^i .

→ We can't bias MD by coords from

ISOMAP / SNE etc.



First major advantage of NNs



in MD we need reduced dim. representation to be explicit & differentiable func² of coords
 \therefore we want to compute forces.