

COMS 4721: Machine Learning for Data Science

Lecture 19, 4/6/2017

Prof. John Paisley

Department of Electrical Engineering
& Data Science Institute
Columbia University

3rd and final matrix factorization technique from a few different perspectives,
but they all relate to something called PCA.

PRINCIPAL COMPONENT ANALYSIS

often used for dimensionality reduction.

DIMENSIONALITY REDUCTION

We're given data x_1, \dots, x_n , where $x \in \mathbb{R}^d$. This data is often

high-dimensional, but the "information" doesn't use the full d dimensions.
Each image as a point in the high dimensional space.



For example, we could represent the above images with three numbers since they have three degrees of freedom. Two for shifts and a third for rotation.

[deciding the center]

Principal component analysis can be thought of as a way of automatically mapping data x_i into some new low-dimensional coordinate system.

- ▶ It captures most of the information in the data in a few dimensions
- ▶ Extensions allow us to handle missing data, and "unwrap" the data.

Comp'n.

* Meaning that vectors in the lower dimensional space all kind of have the same relationship to each other that they had in the higher dimensional space.

And so learning an algorithm on data from the lower dimensional space shouldn't give worse results than what we would get from learning on a higher dimensional space.



how is this possible?

* All we need to learn is a shift and a rotation;
, i.e., 2 shifts and one rotation for these 3 images.

So where is the center, that's a 2-dim vector for the center and then the rotation is a 3rd dimension.

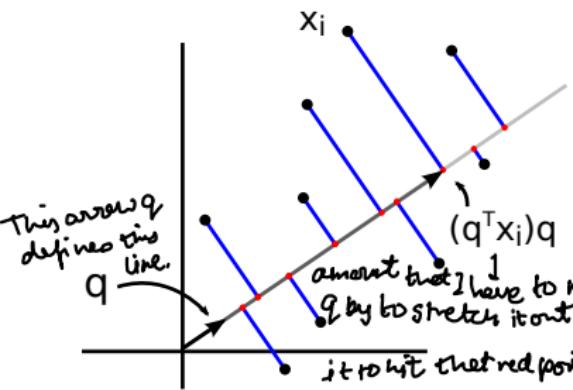
And so we really 3 v.s. for this dataset to capture everything going into it.

And also the base pattern of 3 that's going to be reappearing in each of these images

PRINCIPAL COMPONENT ANALYSIS

Example: How can we approximate this data using a unit-length vector q ?

Represent each of these points in terms of
 q is a unit-length vector, $q^T q = 1$. the vector q is



Red dot: The length, $q^T x_i$, to the axis after projecting x onto the line defined by q .

1 number.

The vector $(q^T x_i)q$ takes q and stretches it by to get the corresponding red dot. (which is the closest we can take q and get it to the point x_i)

So what's a good q ? How about minimizing the squared approximation error,
(a mathematical measure of quality, that we can then minimize to say we found the best vector q)

$$q = \arg \min_q \sum_{i=1}^n \|x_i - qq^T x_i\|^2 \quad \text{subject to } q^T q = 1 \text{ (unit length)}$$

q stretched out by the product $q^T x_i$

$qq^T x_i = (q^T x_i)q$: The approximation of x_i by stretching q to the “red dot.”
P (A problem in a nutshell)

* PCA essentially does this problem.

Let's take this vector Q , make it unit length, and now represent each of these data points by dropping a line perpendicular to the line constructed by the line constructed by the direction of Q .

- ** ↳
- Now we wanna take each of these points and represent it by the amount that we have to stretch Q to get to the red dot by dropping it that line.
 - What does this dot product of two unit vector Q and x_i represent? That represents the length from this origin to the red dot here.

Dimensionality reduction:

So I've done very trivial and not interesting example of dimensionality reduction by taking a two dimensional vector and reducing it now to one dimension.

And so for each of these points, I need to keep only one number. I need to keep the value of $Q^T x_i$.

And then also somewhere, I need to keep the vector Q , potentially if I wanna construct this, try to reconstruct this space.

linear regression
minimizing the sum
of squared errors

dropping along
the o/p dimension.

PCA → dropping
1 to Q

(x_i)

We're taking this two dimensional vector and approximating it with this two dimensional vector. $(Q Q^T x_i)$
But all of those approximations have to fall along the same one dimensional subspace.

same dimensional space
as x_i

PCA : THE FIRST PRINCIPAL COMPONENT

This is related to the problem of finding the largest eigenvalue,

$$q = \arg \min_q \sum_{i=1}^n \|x_i - qq^T x_i\|^2 \quad \text{s.t. } q^T q = 1$$

Draw always doesn't include \$q\$.

now? [vector transposed multiplied by itself.]

constructed by taking each of the our data points \$x_i\$ and putting it along a column

$$= \arg \min_q \sum_{i=1}^n x_i^T x_i - q^T \underbrace{\left(\sum_{i=1}^n x_i x_i^T \right) q}_{d \times n \text{ (each column is a data point)}}.$$

\$d \times n\$ (each column is a data point).

We've defined $X = [x_1, \dots, x_n]$. Since the first term doesn't depend on q and we have a negative sign in front of the second term, equivalently we solve

$$q = \arg \max_q q^T (XX^T)q \quad \text{subject to } q^T q = 1$$

we want to find the vector \$q\$ that maximises this product

This is the eigendecomposition problem: *this product*

► q is the first eigenvector of XX^T *and and symmetric*

► $\lambda = q^T (XX^T)q$ is the first eigenvalue *[the value of this product evaluated at the optimal \$q\$]*

PCA: GENERAL

That's how we would represent every data point with 1 vector. Meaning we take that original data and map it down to 1d space. So we can represent every data point now with 1 no.

The general form of PCA considers K eigenvectors, find K vectors in order to this (approximate x_i by linear combination of K different vectors)

$$q = \arg \min_q \sum_{i=1}^n \|x_i - \underbrace{\sum_{k=1}^K (x_i^T q_k) q_k\|_2^2}_{\text{approximates } x} \quad \text{s.t. } q_k^T q_{k'} = \begin{cases} 1, & k = k' \\ 0, & k \neq k' \end{cases}$$

(unit length)
(* * *)
1 (orthogonal to each other.)

for q *how just adding them up*

$$= \arg \min_q \sum_{i=1}^n x_i^T x_i - \sum_{k=1}^K q_k^T \left(\underbrace{\sum_{i=1}^n x_i x_i^T}_{= XX^T} \right) q_k$$

same as before ↑ *same as before*
but now we have a sum over K different vectors

The vectors in $Q = [q_1, \dots, q_K]$ give us a K dimensional subspace with which to represent the data:

Reduced the no of dimensions to K .

$$x_{\text{proj}} = \begin{bmatrix} q_1^T x \\ \vdots \\ q_K^T x \end{bmatrix},$$

were this over K vectors of subject to 2 constraints
now projection relates to original value
 $x \approx \sum_{k=1}^K (q_k^T x) q_k = Q x_{\text{proj}}$
amount we have to stretch to approximate it.
↳ sum over all vectors.

The eigenvectors of (XX^T) can be learned using built-in software.

* It's the amount that we have to multiply q by to get it as close to the original data as possible.

However, PCA in general doesn't just consider one eigenvector and one eigenvalue.

Meaning that it doesn't map the data to one-dimensional space.
It often will consider K eigenvectors where k is arbitrary

*** We take each eigenvector q_k and we stretch it by the dot product of the point x_i with the eigenvector. So this dot product ($x_i^T q_k$) is the amount that we have to stretch q_k by to get that individual vector as close to x_i as possible.

And now because each of these qs are orthogonal to each other, we treated as K different approximations and then some each of those up. So we're trying to approximate the data with a K dimensional vector, which is the k dimensions that we get from the K different dot products here.

☞ what I gain from moving in 1 dimension has nothing to do with other dimensions.

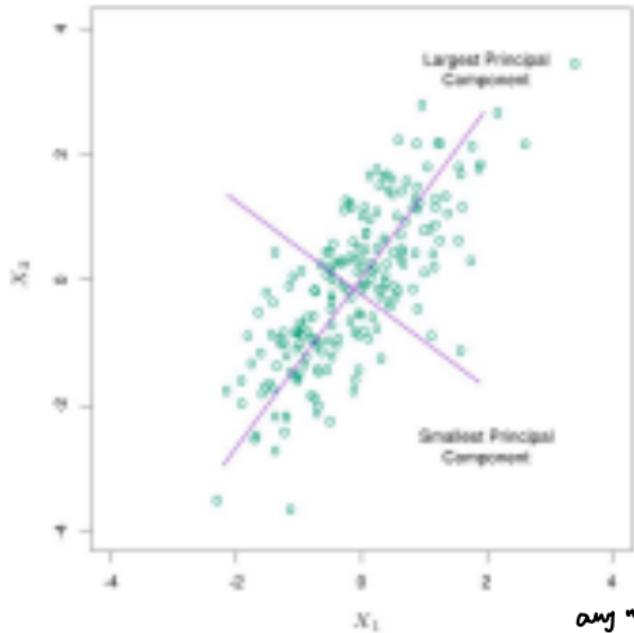
⑤ So now, we've taken the original data and reduce the number of dimensions to K where the projection that we get of the data from the higher dimension is simply represented by taking the dot products of the data with each of these vectors q.

⑥ In practice, when we do this, what we'll often do is we'll take this matrix X, Take a dot product. And then call some function to find the K largest eigenvectors, as well as their associated eigenvalues.

The kth eigenvalue will be the value of this product for q_k^T

$$q_k^T \left(\sum_{i=1}^n x_i x_i^T \right) q_k$$

EIGENVALUES, EIGENVECTORS AND THE SVD



An equivalent formulation of the problem is to find (λ, q) such that

$$(XX^T)q = \lambda q^*$$

Since (XX^T) is a PSD matrix, there are $r \leq \min\{d, n\}$ pairs,

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0,$$

$$q_k^T q_k = 1, \quad q_k^T q_{k'} = 0$$

Why is (XX^T) PSD? Using the SVD, $X = USV^T$, we have that

any matrix $\xrightarrow{\text{columns of } U \text{ are orthonormal}}$ $\xrightarrow{\text{rows of } V^T \text{ are also "}}$ Any

$\xrightarrow{\text{non-negative diagonal values}}$

why?

$$(XX^T) = US^2U^T \Rightarrow Q = U, \quad \lambda_i = (S^2)_{ii} \geq 0 \quad \begin{cases} \text{squared, can't be negative.} \\ \text{eigenvalue} \end{cases}$$

Preprocessing: Usually we first subtract off the mean of each dimension of x . ***

eigen vectors along its columns

- * This matrix (xx^T) here multiplied by λ , it doesn't change its direction. It only stretches it by value λ .

These eigenvalues being positive is a result of this being a positive semidefinite matrix and where we can get are these different eigenvectors where each one is unit length, and are orthogonal to each other.a

- We can find the eigenvectors and the eigenvalues of this matrix product (xx^T) by just taking SVD of the original matrix.

*** Rows of X have their mean subtracted off. So we center the data, so we center the

data so that all of the dimensions in our dataset have mean 0.

↳ Because we want the data to be centered along the original axes. So that we can represent the data in its 1st principal component and its 2nd principal component. So that comes from the fact of these PCs having to pass through the origin

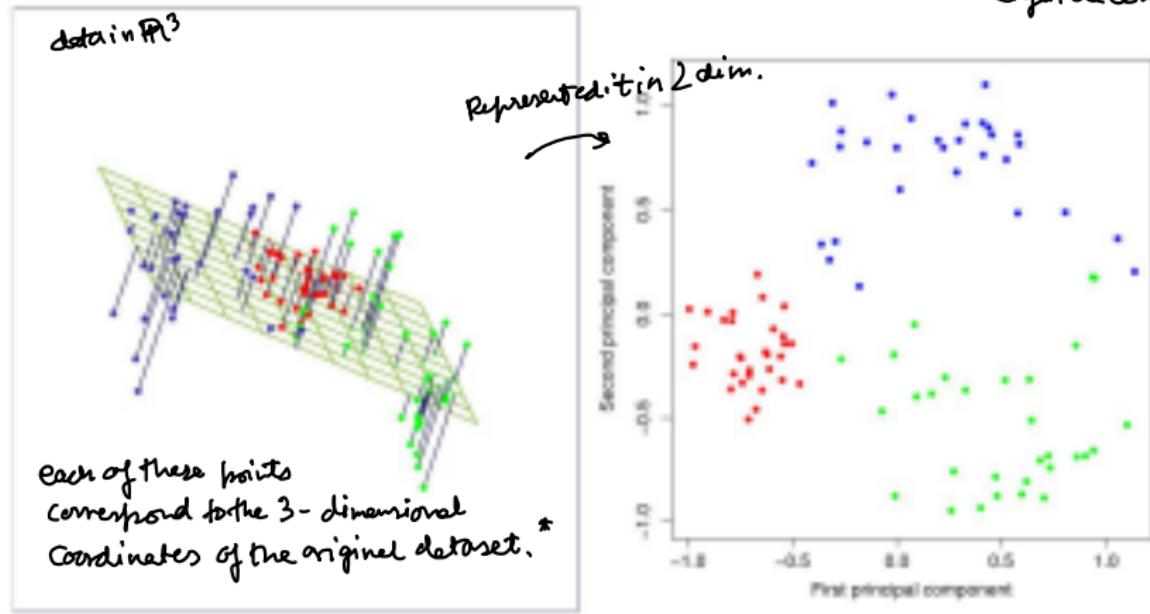
If you imagine this dataset shifted very far away from the origin, it would be much harder to find such a nice representation, if this were rotated a little bit.

[why?]



why
we
do
this?

PCA: EXAMPLE OF PROJECTING FROM \mathbb{R}^3 TO \mathbb{R}^2 using eigen decomposition



For this data, most information (structure in the data) can be captured in \mathbb{R}^2 . **

(left) The original data in \mathbb{R}^3 . The hyperplane is defined by q_1 and q_2 .

(right) The new coordinates for the data: $x_i \rightarrow x_i^{proj} = \begin{bmatrix} x_i^T q_1 \\ x_i^T q_2 \end{bmatrix}$.

* When we try to find an eigen decomposition to project it into two dimensions, it's like finding a plane that slices through the data with a little error as possible.

So, we take a plane. We drop a line to the hyperplane perpendicular to it.

That's our error and we sum that squared error. That's the thing that we're trying to minimize.

**

Notice that if we had just simply taken two of the original dimensions and thrown away the third, we couldn't have gotten as nice of a representation of the dataset.

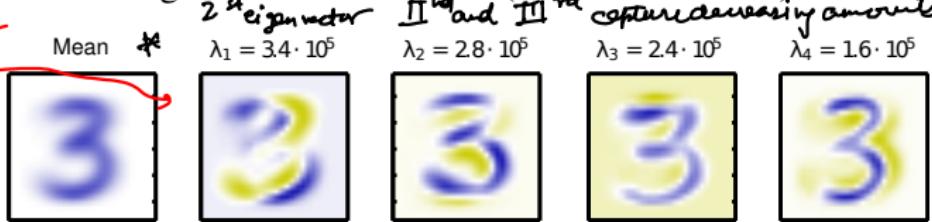
So it's like trying to shift and rotate the data around, so that we can just throw away some of the dimensions but keep the remaining, but where we align  the information as much as possible along the dimensions that we're going to keep  align the information near?

Why are we doing it?

EXAMPLE: DIGITS

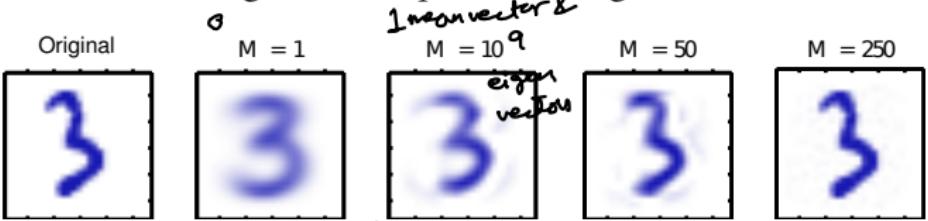
[we don't think 256 nos. are necessarily required to keep all of the information about each of these 3's. So we project them to a lower dimension using eigenvalue decomposition, PCA.] why?

Data: 16×16 images of handwritten 3's (as vectors in \mathbb{R}^{256})



Each eigen vector is 256 dimensions. We roll back into the 16×16 matrix and show it as an image.

Above: The first four eigenvectors q and their eigenvalues λ .



We're almost vertical and the kernel will ignore ones with pick up on the noise here.

Above: Reconstructing a 3 using the first $M - 1$ eigenvectors plus the mean, and approximation

$$x \approx \text{mean} + \sum_{k=1}^{M-1} (x^T q_k) q_k \quad \star \star$$

So this image is a 10 dimensional or really a 9 dimensional representation of the 3. ← how? what is the 9d vector?

* ① So if we took this image and dot multiply that with any other image here and summed it up, it would equal 0.

② And if we took each of these images and dot multiplied it by itself and summed it up, it would equal 1.

Now, we have a 16×16 image, and want to project it into a few dimensions:

1. First, we subtract off the mean from it ← doesn't change the intensity & resulting be not same as original.

2. We take the remaining subtracted image and

dot multiply it by each of these and sum up the results.
(4 eigenvectors)

** We can get now 4 nos. which is the result of dot products of each of these images with the image that just came in.

And now if we want to reconstruct the image, we take the dot product of our means subtracted image, we multiply it by the eigenvector and now that's the approximation like we have discussed previously. of the original image in that dimension only and we sum them up!

o If we wanna represent that only in one vector, we would represent it in the mean

g why? Use
1st eigenvector?

o We are representing this image X as being approximately equal to the mean plus the dot product of X with each eigenvector and we scale the eigenvector and sum it up. This is the approximation we are making.

PROBABILISTIC PCA

We take PCA, and view it as a probability model. And learn the parameters using a probabilistic representation.

PCA AND THE SVD (connection recall)

used the original matrix rather than the dot product.

We've discussed how any matrix X has a singular value decomposition,

$$X = USV^T, \quad U^T U = I, \quad V^T V = I$$

[Assuming X has more columns than rows. If it had more rows than columns then some of these transposes would be switched.]

and S is a diagonal matrix with non-negative entries.

Therefore,
 which want to learn
 eigen decompositions
 \uparrow

$$XX^T = US^2 U^T \Leftrightarrow (XX^T)U = US^2$$

$$(XX^T)q = \lambda q \quad (\text{Representation from previous slide})$$

isn't this true only for PSDs?

U is a matrix of eigenvectors, and S^2 is a diagonal matrix of eigenvalues.
 [eigenvectors along the columns of U .]

[λ along the diag and elements of S^2]

The left singular vectors (u_j) (U) correspond to the principal components of this matrix.
 The square of the singular values correspond to the eigenvalues of this matrix.

A MODELING APPROACH TO PCA

A model that uses the EM algorithm, so we have to define a generative process and then do EM to run a point estimate of it.

Using the SVD perspective of PCA, we can also derive a probabilistic model for the problem and use the EM algorithm to learn it.

This model will have the advantages of: (*that PCA can't handle easily directly*)

- ▶ Handling the problem of missing data
- ▶ Allowing us to learn additional parameters such as noise
- ▶ Provide a framework that could be extended to more complex models
- ▶ Gives distributions used to characterize uncertainty in predictions
- ▶ etc. ↑ Reasons why we might want to choose a probabilistic model.

PROBABILISTIC PCA

(can also view as doing probabilistic SVD because of the connection.)

[What is the significance of V here?]

In effect, this is a new matrix factorization model.

- ▶ With the SVD, we had $X = USV^T$. eigenvectors
eigenvalues [So we can directly learn this sort of factorization without taking the product of X with itself.]
- ▶ Remove all constraints
We now approximate $X \approx WZ$, where what's new
going to do
- depending upon the problem it is being applied to
 - ▶ W is a $d \times K$ matrix. In different settings this is called a "factor loadings" matrix, or a "dictionary." It's like the eigenvectors, but no orthonormality.
 - ▶ The i th column of Z is called $z_i \in \mathbb{R}^K$. Think of it as a low-dimensional representation of x_i . concreteness of value [If the i th column in X corresponds to the i th data point, then the column of Z corresponds to the i th data point except projected into a lower K dimensional space.]

The generative process of Probabilistic PCA is

$$x_i \sim N(W z_i + \text{noise}, \sigma^2 I), \quad z_i \sim N(0, I).$$

In this case, we don't know W or any of the z_i .

Now, we want to learn both of these things.
(We are going to do this with MLE.)

[Multi-variate Gaussian]
the vector z_i is now also a random variable from a Gaussian with zero mean and identity covariance matrix.

? Why this assumption?

• d -dimensionality of data
 k -no. of eigenvalues/eigenvectors we want to learn

• We're also going to relax the constraint that w has to be orthonormal.
So that's a crucial difference between prob. PCA and regular PCA.

The columns don't have to unit length and they don't have to be orthogonal to each other.

• Difference from matrix factorization neither of these have to be non-negative
We want to learn a low rank factorization.

$$X = WZ$$

1000 dimensions
10,000 observations
20
10,000
20
10,000
20

Don't watch!

○ The error in the low-rank approximation is Gram-Schmidt which makes sense since PCA did a squared also (which is like Gram-Schmidt.)

what makes it probably right?

Expectation of $x_i \sim N(Wz_i, \sigma^2 I)$

$\left[\quad \right]$ full &
shiny

$z_i \sim N(0, I)$

as a prior

THE LIKELIHOOD

Maximum likelihood

Our goal is to find the maximum likelihood solution of the matrix W under the marginal distribution, i.e., with the z_i vectors integrated out, ← why?

Then do a point estimate that maximizes (integrate out impact of z_i)
the marginal likelihood with respect to matrix W .

$$W_{\text{ML}} = \arg \max_W \ln p(x_1, \dots, x_n | W) = \arg \max_W \sum_{i=1}^n \ln p(x_i | W). \quad \text{①}$$

product → sum
all of the data vectors.

independence assumption

?

This is intractable because $p(x_i | W) = N(x_i | 0, \underbrace{\sigma^2 I + WW^T}_{\text{covariance matrix}}$,
density:

$$N(x_i | 0, \sigma^2 I + WW^T) = \frac{1}{(2\pi)^{\frac{d}{2}} |\sigma^2 I + WW^T|^{\frac{1}{2}}} e^{-\frac{1}{2}x^T(\sigma^2 I + WW^T)^{-1}x}$$

We can set up an EM algorithm that uses the vectors z_1, \dots, z_n .

Even though we couldn't take the derivative of this with respect to W & optimize it, we can still evaluate it at each value of W . We can use this eqn. ① where we take the log of this ① over each.

reintroduce these vectors as
missing part of our model.
and then do EM on that.

If we take the log of it, we
can't do this and solve for
 W analytically.

EM FOR PROBABILISTIC PCA

Setup up the EM equality

The marginal log likelihood can be expressed using EM as

$$\sum_{i=1}^n \ln \int p(x_i, z_i | W) dz_i \xrightarrow{\substack{\text{marginal likelihood which is} \\ \text{hard to optimise over } W}} \sum_{i=1}^n \int q(z_i) \ln \frac{p(x_i, z_i | W)}{q(z_i)} dz_i \xleftarrow{\substack{\text{introduced a } q \text{ distribution over each of these} \\ \text{vectors } z_i; \\ \text{joint likelihood of } x_i \text{ & } z_i}} \mathcal{L}$$

Marginal likelihood we want
to minimize over W as a
integral $\int \dots dz_i$ is integrated
out. [Now if we want to use the RHS in order to minimize the
LHS over W .]

$$+ \sum_{i=1}^n \int q(z_i) \ln \frac{q(z_i)}{p(z_i | x_i, W)} dz_i \xleftarrow{\text{KL}}$$

divergence of q distribution
on z_i and conditional posterior
distribution of z_i given
the matrix W .

EM Algorithm: Remember that EM has two iterated steps

- E-step
- Set $q(z_i) = p(z_i | x_i, W)$ for each i (making KL = 0) and calculate \mathcal{L}
 - Maximize \mathcal{L} with respect to W
- M-step

Again, for this to work well we need that

in closed form.

- we can calculate the ^{conditional} posterior distribution $p(z_i | x_i, W)$, and
- maximizing \mathcal{L} is easy, i.e., we update W using a simple equation

THE ALGORITHM

EM for Probabilistic PCA

Given: Data $x_{1:n}, x_i \in \mathbb{R}^d$ and model $x_i \sim N(Wz_i, \sigma^2 I)$, $z_i \sim N(0, I)$, $z \in \mathbb{R}^K$

Output: Point estimate of W and posterior distribution on each z_i about what the likelihood which gives us our posterior belief of what the low-dimensional embedding should be.

E-Step: Set each $q(z_i) = p(z_i | x_i, W) = N(z_i | \mu_i, \Sigma_i)$ where $\mu_i = \Sigma_i W^T x_i / \sigma^2$ and $\Sigma_i = (I + W^T W / \sigma^2)^{-1}$.

every single z_i has exactly the same covariance in the conditional posterior noise parameter we can show it's a multi-variate gaussian

M-Step: Update W by maximizing the objective \mathcal{L} from the E-step calculate using the next recent value of W in both of these equations.

L get this closed form update of $W = \left[\sum_{i=1}^n x_i \mu_i^T \right] \left[\sigma^2 I + \sum_{i=1}^n (\mu_i \mu_i^T + \Sigma_i) \right]^{-1}$ we update W after updating each of these $q(z_i)$ s.

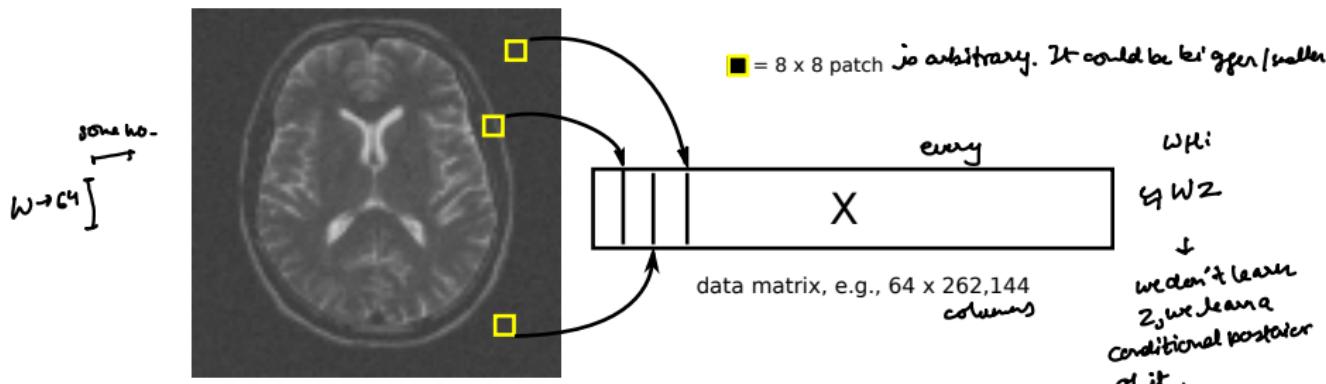
Iterate E and M steps until increase in $\sum_{i=1}^n \ln p(x_i | W)$ is “small.”

Comment:

- ▶ The probabilistic framework gives a way to learn K and σ^2 as well.

* So we don't just get a low embedding for each point x_i ; in the \mathcal{G} of z_i , we also get some distribution on that embedding.

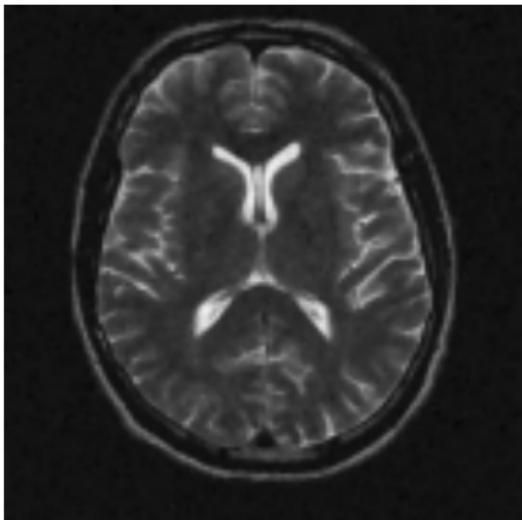
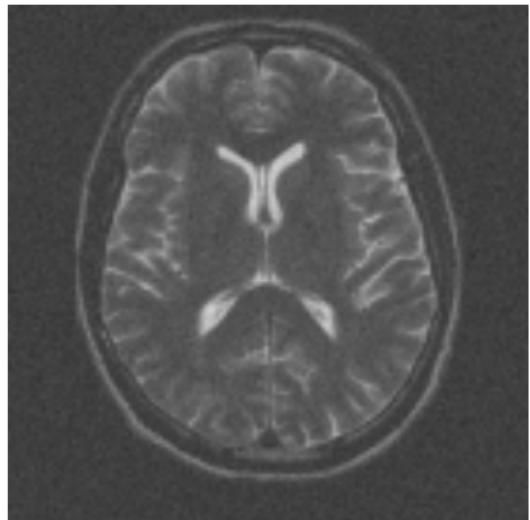
EXAMPLE: IMAGE PROCESSING problem of denoising an image.



For image problems such as denoising or inpainting (missing data)

- ▶ Extract overlapping patches (e.g., 8×8) and vectorize to construct X
Factorise using
- ▶ Model with a factor model such as Probabilistic PCA
- ▶ Approximate $x_i \approx W\mu_i$, where μ_i is the posterior mean of z_i
^{n conditional posterior mean}
how?
To do this, just plug in the
mean for all the
uncertainty we had.
- ▶ Reconstruct the image by replacing x_i with $W\mu_i$ (and averaging)
Noise is gone in noise term of the
generative model.

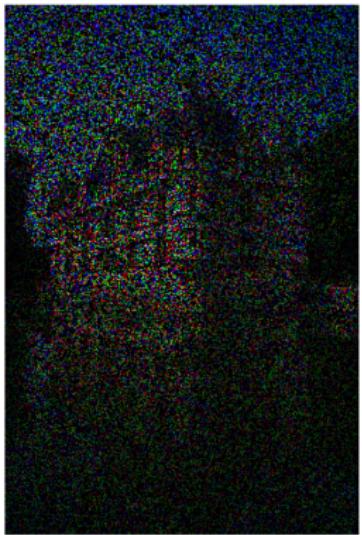
EXAMPLE: DENOISING



Noisy image on left, denoised image on right. The noise variance parameter σ^2 was learned for this example.

EXAMPLE: MISSING DATA

by throwing away 80%. If so, you didn't throw away any information because all information is in a low-dimensional image.



Another somewhat extreme example:

- ▶ Image is $480 \times 320 \times 3$ (RGB dimension)
- ▶ Throw away 80% ^{of value} at random
- ▶ (left) Missing data, (middle) reconstruction, (right) original image

Extract a $8 \times 8 \times 3$ cube and
refill in the missing ~~area~~ ?
data using EM.

KERNEL PCA

KERNEL PCA

We've seen how we can take an algorithm that uses dot products, $x^T x$, and generalize with a nonlinear kernel. This generalization can be made to PCA.

replace dot products with $\langle \cdot, \cdot \rangle$

$x = \text{each } x_i \text{ along its column dimension}$
 $\text{to higher dimensions}$ \rightarrow XX^T \rightarrow eigenvectors
 outer product \rightarrow of this matrix .

Recall: With PCA we find the eigenvectors of the matrix $\sum_{i=1}^n x_i x_i^T = XX^T$ *and the eigenvectors of this matrix.*

- ▶ Let $\phi(x)$ be a feature mapping from \mathbb{R}^d to \mathbb{R}^D , where $D \gg d$
Projecting to much higher dimensions before projecting to much lower dimension can do even more for us.
- ▶ We want to solve the eigendecomposition

$$\left[\sum_{i=1}^n \phi(x_i) \phi(x_i)^T \right] q_k = \lambda_k q_k$$

$(XX^T)q_k = \lambda_k q_k$

Find the vector q and λ for matrix from this higher dimensional projection.

without having to work in the higher dimensional space.

(For e.g. in Gaussian kernels, this is an ∞ -dim space.)

- ▶ That is, how can we do PCA without explicitly using $\phi(\cdot)$ and q ?

KERNEL PCA

Notice that we can reorganize the operations of the eigendecomposition

$$k \xrightarrow{i \text{ indexing}} \text{data point} \quad k \xrightarrow{\text{eigenvalue, eigenvector pair.}} \sum_{i=1}^n \phi(x_i) \underbrace{(\phi(x_i)^T q_k) / \lambda_k}_{= a_{ki}} = q_k$$

λ_k on the left side by dividing it.

Multiply higher dimensional projection of x_i with k^{th} eigenvector, divided by k^{th} eigenvalue.

That is, the eigenvector $q_k = \sum_{i=1}^n a_{ki} \phi(x_i)$ for some vector $a_k \in \mathbb{R}^n$.
 Represent each eigenvector as a linear combination of data projected to the higher dimensional space.

The trick is that instead of learning q_k , we'll learn a_k .

Plug this equation for q_k back into the first equation:

(Why are you plugging it back into the same equation?)

$$\Phi^T \sum_{i=1}^N \phi(x_i) \sum_{j=1}^n a_{kj} \underbrace{\phi(x_i)^T \phi(x_j)}_{= K(x_i, x_j)} = \Phi^T \lambda_k \sum_{i=1}^n a_{ki} \phi(x_i)$$

dots products! (back)

Φ vector in $\mathbb{R}^{n \times n}$ $n \rightarrow$ no. of data points

Trick: Instead of learning this potentially infinite dimensional vector, we're going to learn an n -dimensional vector corresponding to that k^{th} eigenvector.

and multiply both sides by $\phi(x_l)^T$ for each $l \in \{1, \dots, n\}$.

$$\Phi = [\phi(x_1) \ \phi(x_2) \ \dots \ \phi(x_n)]$$

projection of the l^{th} data point, for each of the n data points.

KERNEL PCA

(will simplify our problem)

When we multiply the following by $\phi(x_l)^T$ for $l = 1 \dots, n$:

$$\Phi = [\phi(x_1) \dots \phi(x_n)]$$

$$\underbrace{\Phi^T}_{\Phi = [\phi(x_1) \dots \phi(x_n)]} \sum_{i=1}^N \phi(x_i) \sum_{j=1}^n a_{kj} \underbrace{\phi(x_i)^T \phi(x_j)}_{= K(x_i, x_j)} = \underbrace{\Phi^T}_{\Phi^T} \lambda_k \sum_{i=1}^n a_{ki} \phi(x_i)$$

we get a new set of linear equations

multiply both sides by K^{-1} (PSD matrix)

$$K^2 \mathbf{a}_k = \lambda_k K \mathbf{a}_k \iff K \mathbf{a}_k = \lambda_k \mathbf{a}_k \quad \begin{matrix} \text{Kernel} \\ K_{ij} \rightarrow \text{data point } x_i \cdot x_j \end{matrix}$$

where K is the $n \times n$ kernel matrix constructed on the data.

an eigenvector associated with the kernel matrix constructed among data

Because K is guaranteed to be PSD because it is a matrix of dot-products, the LHS and RHS above share a solution for $(\lambda_k, \mathbf{a}_k)$.

Now perform “regular” PCA, but on the kernel matrix K instead of the data matrix XX^T . We summarize the algorithm on the following slide.

originally, I could have mapped the data into such a high dimensional space or infinite dimensional space

KERNEL PCA ALGORITHM

Kernel PCA

Given: Data $x_1, \dots, x_n, x \in \mathbb{R}^d$, and a kernel function $K(x_i, x_j)$.

define
all points.

Construct: The kernel matrix on the data, e.g., $K_{ij} = b \exp\left\{-\frac{\|x_i - x_j\|^2}{c}\right\}$.

Solve: The eigendecomposition

$$Ka_k = \lambda_k a_k$$

dimensionality we want to project the data into.

for the first $r \ll n$ eigenvector/eigenvalue pairs $(\lambda_1, a_1), \dots, (\lambda_r, a_r)$. of kernel K

Output: A new coordinate system for x_i by (implicitly) mapping $\phi(x_i)$ and then projecting $q_k^T \phi(x_i)$

back to the
lower dimensional
space -
 x_i projection
(r -dimensional
space.)

$\lambda_1 a_{1i}$
 \vdots
 $\lambda_r a_{ri}$

r different eigen vectors there
For the k^{th} dimension of our projection,
we learnt the k^{th} eigenvector or
eigenvalue.
For the i^{th} datapoint we look at the
 i^{th} dimension of k^{th} eigenvector and
multiply it.

where a_{ki} is the i^{th} dimension of the k^{th} eigenvector a_k .

KERNEL PCA AND NEW DATA

→ how to project it using the same assumptions as the model was.

Q: How do we handle new data, x_0 ? Before, we could take the eigenvectors q_k and project $x_0^T q_k$, but a_k is different here.

that's the k^{th} dimension of the lower dimensional projection

A: Recall the relationship of a_k to q_k in kernel PCA is

k^{th} eigenvector in the higher dimensional space can be written as:

$$q_k = \sum_{i=1}^n a_{ki} \phi(x_i). \quad \begin{matrix} \xrightarrow{k \rightarrow \text{identity eigenvector}} \\ \downarrow i \rightarrow \text{"data point"} \end{matrix}$$

We used the “kernel trick” to avoid working with or even defining $\phi(x_i)$.

As with regular PCA, after mapping x_0 we want to project onto eigenvectors

$$x_0 \xrightarrow{\text{projection}} \begin{bmatrix} \phi(x_0)^T q_1 \\ \vdots \\ \phi(x_0)^T q_r \end{bmatrix} \quad \begin{matrix} \xrightarrow{\text{1st dimension in the lower}} \\ \text{dimensional space.} \end{matrix}$$

Plug from ①

k^{th} eigenvector of the kernel matrix constructed on n data points.

Plugging in for q_k : $\phi(x_0)^T q_k = \sum_{i=1}^n a_{ki} \phi(x_0)^T \phi(x_i) = \sum_{i=1}^n a_{ki} K(x_0, x_i).$

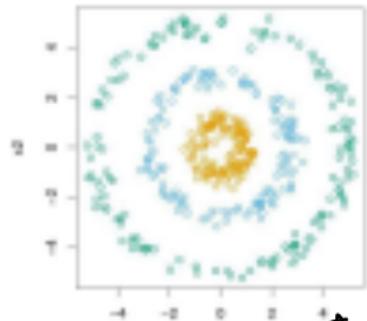
Kernel to evaluate using new data point and all old data points in my dataset.

* So I evaluate the kernel between a new incoming point with all the points in my dataset.
I multiply the kernel by the eigenvector I got from the matrix, (the original matrix K),
sum it up.

That's my projection.

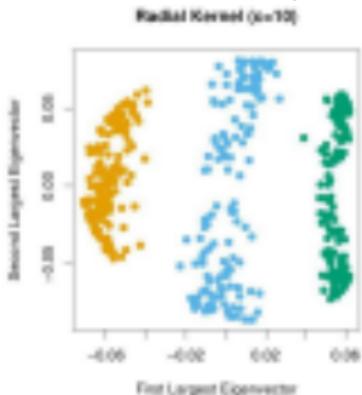
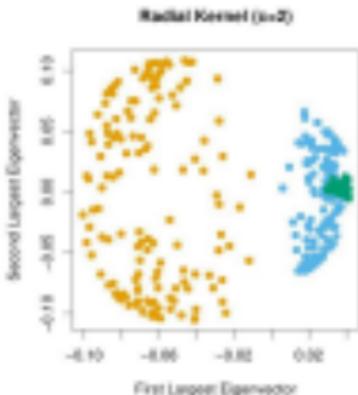
EXAMPLE RESULTS

color coded just for visualization purposes



means here wouldn't
seem 3 ring

Project into 2 dimensions : there would be no projection.
If we want to separate these 3 rings from each other : take these data
points, and project it into a higher dimensional space, and then do
PCA there .



mean world here

An example of kernel PCA using the Gaussian kernel.

- (left) Original data, colored for reference (but may be classes)
- (middle) New coordinates using kernel width $c = 2$
- (right) New coordinates using kernel width $c = 10$

Terminology: What we are doing is closely related to “spectral clustering”
and can be considered an instance of “manifold learning.”

Take each of these data points and calculating the kernel matrix. The pairwise kernel matrix
describing a Gaussian kernel and then using that $n \times n$ kernel matrix, projecting onto its
 Σ^{+2} eigenvectors and eigenvalues.