ICA = PP

ICA = PP

M.C. JONES and ROBIN SIBSON. What is projection Pursuit? Journal of the Royal Statistical Society. 1987

Aapo Hyvärinen. Fast and Robust Fixed-Point Algorithms for Independent Component Analysis. IEEE TRANSACTIONS ON NEURAL NETWORKS, 1999

How to make sense of a table of data?



Jonathan Power "The Plot"

What do I do with this data?



91282 columns

Take the mean?



91282 columns



Take the mean?





What have we done?



What have we done?



= w^T X Orthogonal projection of X onto w

What have we done?





:





How do we find w?

 Solution: project onto all the direction and look at the histogram for some interesting features Plan of talk:

- Defining "interestingness"
- PP
- ICA

Notation

- Data is a matrix X (dimensions: features x samples)
- Direction vector is w (features x 1)
- If we want more that one vector then it is W (features x p)







$$\mathbf{X}_c = \mathbf{X}(\mathbf{I} - \mathbf{1}\mathbf{1}^T/\mathbf{n})$$





Location of data should not matter --> centering

$$\mathbf{X}_c = \mathbf{X}(\mathbf{I} - \mathbf{1}\mathbf{1}^T/\mathbf{n})$$

Orientation of data should not matter --> Ok.
 (Rw)^TRX





Location of data should not matter --> centering

$$\mathrm{X}_{c} = \mathrm{X}(\mathbf{I} - \mathbf{1}\mathbf{1}^{T}/\mathrm{n})$$

Orientation of data should not matter --> Ok.
 (Rw)^TRX



• Variance?

max Variance($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$) = E[$(\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c})^{2}$]

max Variance($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$) = E[($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$)²]

$$= 1/n^* (\mathbf{w}^\mathsf{T} \mathbf{X}_c) (\mathbf{w}^\mathsf{T} \mathbf{X}_c)^\mathsf{T}$$

- max Variance($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$) = E[$(\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c})^{2}$]
 - $= 1/n^* (\mathbf{w}^\mathsf{T} \mathbf{X}_c) (\mathbf{w}^\mathsf{T} \mathbf{X}_c)^\mathsf{T}$

$$= 1/n^* \mathbf{w}^{\mathsf{T}} (\mathbf{X}_{c} \mathbf{X}_{c}^{\mathsf{T}}) \mathbf{w}$$

max Variance($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$) = E[($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$)²]

$$= 1/n^* (\mathbf{w}^\mathsf{T} \mathbf{X}_c) (\mathbf{w}^\mathsf{T} \mathbf{X}_c)^\mathsf{T}$$

=
$$1/n^* w^T (X_c X_c^T) w$$

 \downarrow
1st eigenvector of covariance matrix

PCA!

max Variance(
$$\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$$
) = E[($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$)²]

$$= 1/n^* (\mathbf{w}^\mathsf{T} \mathbf{X}_c) (\mathbf{w}^\mathsf{T} \mathbf{X}_c)^\mathsf{T}$$

=
$$1/n^* \mathbf{w}^T (\mathbf{X}_c \mathbf{X}_c^T) \mathbf{w}$$

1st eigenvector of covariance matrix

PCA!

(theorem: a symmetric matrix M [like XX^T] has all the eigenvectors and one of them points in the direction of maximum x^TMX)

- Capture largest variance
- Incidentally, approximate the data best



- Capture largest variance
- Incidentally, approximate the data best

$w^T X$ 1D version of X



- Capture largest variance
- Incidentally, approximate the data best

- w^TX 1D version of X
- ww^TX same dimensions as X



- Capture largest variance
- Incidentally, approximate the data best

- w^TX 1D version of X
- ww^TX same dimensions as X
- $|X-ww^TX|^2$ minimize this wrt w?



- Capture largest variance
- Incidentally, approximate the data best

- w^TX 1D version of X
- ww^TX same dimensions as X
- $|X-ww^TX|^2$ minimize this wrt w?

```
trace( (X-ww^TX)(X-ww^TX)^T)
```



- Capture largest variance
- Incidentally, approximate the data best

- w^TX 1D version of X
- ww^TX same dimensions as X
- $|X-ww^TX|^2$ minimize this wrt w?

trace($(X-ww^TX)(X-ww^TX)^T$) trace($-ww^TXX^T$ $-XX^Tww^T$ $+ww^TXX^Tww^T$)

- Capture largest variance
- Incidentally, approximate the data best

- w^TX 1D version of X
- ww^TX same dimensions as X
- $|X-ww^TX|^2$ minimize this wrt w?

trace($(X-ww^TX)(X-ww^TX)^T$) trace($-ww^TXX^T$ $-XX^Tww^T$ $+ww^TXX^Tww^T$) trace($-w^TXX^Tw$)



PCA interpretations

- Capture largest variance
- Best approximation to data
- Direction that is most aligned with data



PCA interpretations

- Capture largest variance
- Best approximation to data
- Direction that is most aligned with data

w^TXX^Tw = sum_squares(X^Tw)



The problem with PCA

The direction of maximum variance is not always interesting



Interestingness

 $I(w^T X)$

Location of data should not matter --> centering

 $\mathbf{X}_c = \mathbf{X}(\mathbf{I} - \mathbf{1}\mathbf{1}^T/\mathbf{n})$



- Orientation of data should not matter --> Ok.
 (Rw)^TRX
- Variance?

Sphering

n.k.a. whitening

• How to make it so

Variance(**w^TX**_c) = 1

for all **w** ?

Sphering

n.k.a. whitening

• How to make it so

Variance($\mathbf{w}^{\mathsf{T}}\mathbf{X}_{c}$) = 1 for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$
n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

 $\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}$ eigenvalue decomposition

n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

 $\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}$ eigenvalue decomposition

 $w^{T}(QX_{c}X_{c}^{T}Q^{T})w = w^{T}QUDU^{T}Q^{T}w$

n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

 $\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}$ eigenvalue decomposition

 $w^{\mathsf{T}}(\mathbf{Q}\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}})w = w^{\mathsf{T}}\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}}w$ $\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}} = \mathsf{Identity}$

n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

 $\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}$ eigenvalue decomposition

 $w^{\mathsf{T}}(\mathbf{Q}\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}})w = w^{\mathsf{T}}\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}}w$ $\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}} = \text{Identity}$

 $Q = UD^{-1/2}U^{T}$

n.k.a. whitening

• How to make it so $Variance(\mathbf{w}^T \mathbf{X}_c) = 1$ for all \mathbf{w} ?

 $X_c \rightarrow X_s = QX_c$

n* Variance(w^TQX_c) = $w^T(QX_c(QX_c)^T)w = w^TQX_cX_c^TQ^Tw$

 $\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}$ eigenvalue decomposition

$w^{\mathsf{T}}(\mathbf{Q}\mathbf{X}_{c}\mathbf{X}_{c}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}})w = w^{\mathsf{T}}\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}}w$ $\mathbf{Q}\mathbf{U}\mathbf{D}\mathbf{U}^{\mathsf{T}}\mathbf{Q}^{\mathsf{T}} = \text{Identity}$

Q=UD-1/2**U**T project onto eigenvectors rescale (undo what X does) project back

Centering



function X_c = center(X)

$$n = size(X,2);$$

 $C = eye(n) - ones(n,1)*ones(1,n)/n;$
 $X_c = X*C;$

function $X_c = center(X)$ $X_c = X - mean(X,2);$

Sphering whitening



function X_s = whiten(X)

[V,D] = eig(cov(X')); Q = V*diag(1./sqrt(diag(D)))*V'; X_s = Q*center(X);

centered

whitened



variance long different orientations

function plot_var(X)

```
angles = linspace(0,2*pi,1000); % angles to calculate variance along
V = zeros(size(angles)); % Variance
for i=1:length(angles)
 % get vector w along angle
 [w_x,w_y] = pol2cart(angles(i),1);
 w = [w_x;w_y];
 % calculate cf along w
 V(i) = var(w'*X);
end
figure,hold on
 plot(X(1,:),X(2,:),'bo')
```

[xx,yy] = pol2cart(angles,V); plot(xx,yy,'linewidth',2,'color','k','linestyle','--') axis equal grid on

Interestingness

 $I(w^T X)$

What is uninteresting?

Gaussian I am afraid

(everything is Gaussian by default)

 If you randomly project X, chances are you get a Gaussian (central limit theorem)

- If you randomly project X, chances are you get a Gaussian (central limit theorem)
- Gaussian is the maximum entropy distribution (if you fix the mean and variance)

- If you randomly project X, chances are you get a Gaussian (central limit theorem)
- Gaussian is the maximum entropy distribution (if you fix the mean and variance)
- So.... minimize entropy! (maximize neg-entropy). I.e. maximize non-gaussianity.

Relative entropy = $-\int p \log(q/p)$

Distance between p and q (>=0)

Relative entropy = $-\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

Relative entropy $= -\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

Relative entropy $= -\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q \log(q) - Entropy(p)$

Relative entropy $= -\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q \log(q) - Entropy(p)$

q Gaussian -> log(q) is quadratic

Relative entropy $= -\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q \log(q) - Entropy(p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

q Gaussian -> log(q) is quadratic -> $E_p[log(q)] = E_q[log(q)] = combination of$ 1st and 2nd moments!

Relative entropy $= -\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q \log(q) - Entropy(p)$

= Entropy(q) - Entropy(p)

q Gaussian -> log(q) is quadratic

 $-> E_p[log(q)] = E_q[log(q)] = combination of 1st and 2nd moments!$

Relative entropy = $-\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q.log(q) - Entropy(p)$

= Entropy(q) - Entropy(p)

>=0

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

q Gaussian -> log(q) is quadratic -> $E_p[log(q)] = E_q[log(q)] = combination of$ 1st and 2nd moments!

Relative entropy = $-\int p \log(q/p)$

Entropy = $-\int p \log(1/p)$

Distance between p and q (>=0)

A bit like distance between p and 1 (no structure)

 $-\int p \log(q/p) = -\int p \log(q) -\int p \log(1/p)$

 $= -\int q \log(q) - Entropy(p)$

= Entropy(q) - Entropy(p)

>=0

q Gaussian -> log(q) is quadratic -> $E_p[log(q)] = E_q[log(q)] = combination of$ 1st and 2nd moments!

Gaussian = maximum entropy when mean and variance fixed

hmm... we are given samples, not the actual data distribution

- hmm... we are given samples, not the actual data distribution
- kernel density not an option (we need to find w with optimization)

- hmm... we are given samples, not the actual data distribution
- kernel density not an option (we need to find w with optimization)
- approximation

 $\mathbf{I}(\mathbf{w}^T \mathbf{X}) = E[(\mathbf{w}^T \mathbf{X})^2]$

variance

 $I(\mathbf{w}^T \mathbf{X}) = E[g(\mathbf{w}^T \mathbf{X})]$

other (as long as g is not quadratic)

Advantages of approximation: continuous function (good for optimization) easy to calculate (unlike entropy)



Let's try this

E.g. rFMRI spatial ICA: each point = voxel each axis = time point



Random projection w'X









Gauss



pow3 (skewness)



logcosh (tanh)



pow4 (kurtosis)



Beyond 2D

- We can't just look at the 2D projections
- We need an algorithm fo finding w
- We have a cost function! $Eig[g(\mathbf{w}^T\mathbf{X})ig]$ s.t. $\mathbf{w}^T\mathbf{w}=1$
- FastICA = Newton-like algorithm

ICA

data mixed signals **X = AS** mixing matrix
ICA











 $W = A^{-1}$



w X = S



Temporal PP

S (2xn)

AS







Successful unmixing

Still some mixing







Interpretation: similarity between voxel time course and IC time course





• When we talk about IC's what are we talking about?

- When we talk about IC's what are we talking about?
- Remember: X=AS

- When we talk about IC's what are we talking about?
- Remember: X=AS
- but X is sphered: QX=AS

- When we talk about IC's what are we talking about?
- Remember: X=AS
- but X is sphered: QX=AS
- $A = mixing matrix = W^{-1}$

- When we talk about IC's what are we talking about?
- Remember: X=AS
- but X is sphered: QX=AS
- $A = mixing matrix = W^{-1}$
- --> X=Q⁻¹AS

- When we talk about IC's what are we talking about?
- Remember: X=AS
- but X is sphered: QX=AS
- $A = mixing matrix = W^{-1}$
- --> X=Q⁻¹AS
- rows of S (independent components) are actually orthogonal!

- When we talk about IC's what are we talking about?
- Remember: X=AS
- but X is sphered: QX=AS
- $A = mixing matrix = W^{-1}$
- --> X=Q⁻¹AS
- rows of S (independent components) are actually orthogonal!
- This is because: SS^T=(WQX)^TWQX=WQXX^TQ^TW^T=WW^T=Identity

 ICA (like PCA) has an explicit model: X = A*S whereas PP only has data and tries to find w s.t. w'X has a funky distribution

- ICA (like PCA) has an explicit model: X = A*S whereas PP only has data and tries to find w s.t. w'X has a funky distribution
- Because ICA has a model of the data, it can be augmented with e.g. noise model (like Melodic), explicit model for S (e.g. mixture of gaussians like FLICA or other mixtures like PROFUMO), etc.

- ICA (like PCA) has an explicit model: X = A*S whereas PP only has data and tries to find w s.t. w'X has a funky distribution
- Because ICA has a model of the data, it can be augmented with e.g. noise model (like Melodic), explicit model for S (e.g. mixture of gaussians like FLICA or other mixtures like PROFUMO), etc.
- But FastICA is PP with nice approximations to negentropy and an efficient algorithm for finding w







The end.