

Principal component analysis

Ferath Kherif, Adeliya Latypova

Laboratory of Research in Neuroimaging (LREN), Department of Clinical Neurosciences, CHUV, University of Lausanne, Lausanne, Switzerland

12.1 Introduction

Principal component analysis (PCA) is the most used method for data exploration and data analysis across all fields of science (Jolliffe, 1986). PCA belongs to the family of dimension reduction methods and is particularly useful when the data at hand are large (i.e., multiple variables), big (i.e., multiple observations per variable), and highly correlated. With such high-dimensional data, the goal is to identify a reduced set of features that represent the original data in a lower-dimensional subspace with a minimal loss of information (Fig. 12.1). There are multiple advantages in dealing with a reduced dataset instead of the original high-dimensional data. These advantages include the following:

- The ability to visualize the data in 2D or 3D
- Less storage space

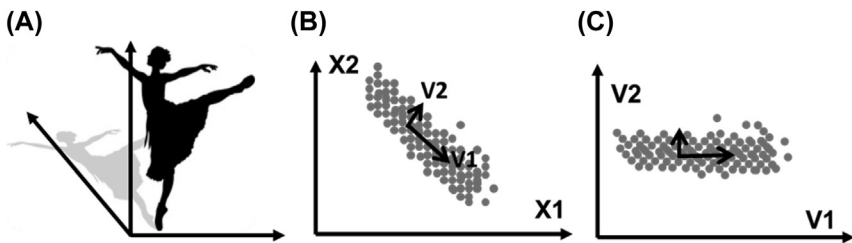


FIGURE 12.1 Principal component analysis (PCA) projections: the aim of the PCA is to find orthogonal projections that can explain the best the variance of the data (A). V1, the first principal axis explains most of the variance of data; V2 is the second principal axis and explains the remaining variance not explained by V1 (B). V1 and V2 form a new orthonormal coordinate system where the data can be projected into (C).

- Removal of collinearity
- Noise reduction

As a dimension reduction method, PCA projects the data into a new, lower-dimensional subspace. Original observations of possibly correlated variables are transformed into a set of values of linearly uncorrelated variables called principal components. Specifically, the original data are projected into a new coordinate system where the first axis, called the first principal axis, corresponds to the direction along which the data vary the most; the second axis, called the second principal axis, corresponds to the direction along which the data vary the most after the first direction; etc. The first principal component is the projection of the original data to the first principal axis and captures the greatest amount of the variance in the data. The second principal component is the projection of the original data to the second principal axis and explains the greatest amount of the variance in the data that is not captured by the first principal component. Each subsequent principal component explains the greatest amount of variance possible under the constraint that it is orthogonal to the preceding principal components, until all the data matrix is decomposed. The total number of principal components corresponds to the number of dimensions in the original data. Critically, because the last principal components will capture the lowest variance in the data, they can be omitted. This is how the lower dimensionality is reached. In the literature, the term principal components and principal axes are often used interchangeably.

Formally, the principal components are obtained by first deconstructing the original data into eigenvectors and eigenvalues. An eigenvector corresponds to the direction with the greatest variance in the data (confusingly, the same term is often used interchangeably with the term principal component). For example, the eigenvectors of [Fig. 12.1A](#) can be represented as V1 and V2 in [Fig. 12.1B](#). Each eigenvector has a corresponding eigenvalue. An eigenvalue is a number that indicates the amount of variance in the data along its corresponding eigenvector. Therefore, the eigenvector with the highest eigenvalue will be the first principal component; the eigenvector with the second highest eigenvalue will be the second principal component; etc. Once estimated, the principal components are used to create a new coordinate space where the data can be projected into ([Fig. 12.1C](#)).

Although PCA is a linear method—it only identifies new dimensions that are the linear combinations of the original ones—it remains one of the most powerful techniques for dimension reduction. The fact that these new dimensions are simple to compute makes PCA the method of choice not only for data exploration but also for data preprocessing before applying more complex statistical or machine learning tools. Dimension reduction via PCA allows to:

$$\begin{array}{c} Y \\ (n \times p) \end{array} = \lambda_1 \begin{array}{c} U(:,1) \end{array} * \begin{array}{c} V(1,:) \end{array} + \lambda_2 \begin{array}{c} U(:,2) \end{array} * \begin{array}{c} V(2,:) \end{array} + \dots + \lambda_p \begin{array}{c} U(:,p) \end{array} * \begin{array}{c} V(p,:) \end{array}$$

FIGURE 12.2 Rank-1 decomposition using singular value decomposition. The data matrix Y can be decomposed in principal components and directions.

- Obtain a better insight into the data which could help generate new hypotheses
- Identify potential issues with the data such as artifacts or outliers
- Reduce the number of predictors of a linear regression model, therefore avoiding the multicollinearity problem and reducing the risk of overfitting
- Obtain higher accuracy and efficiency for all the methods based on computing distances between data points (e.g., K -nearest neighbors [KNN], K -means, Support Vector Machine [SVM])
- Enable the algorithm applied to the data to run faster, which provides additional time for parameters optimization or models benchmarking

The increasing collection of comprehensive “omic” data (genome, proteome, metabolome, microbiome, neuroimaging) provides the opportunity to investigate the biological mechanisms of brain disorders in greater detail than ever before. However, to develop clinically useful models of psychiatric and neurological disease, we need to integrate “omic” data with clinical measures of the manifestation and progression of the illness. In addition, in light of accumulating evidence that biological, psychological, and social factors interact throughout the course of life, we must integrate these factors for a comprehensive understanding of the manifestation and progression of disease. PCA can play an important role for developing disease models that meet these criteria. More importantly, in the era of Big Data and increasing interest in personalized medicine, PCA has the unique advantage of providing a compact representation of the data that captures the main components of individual differences. In the following section, we first explain the mathematical formula behind PCA. In the second section, we show how to implement the method using a toy example with a dementia dataset. Finally, in the last section, we discuss some exemplar applications of PCA to the investigation of brain disorders from the existing literature.

12.2 Method description

The computation of the principal components is relatively simple. It relies on basic principles of linear algebra such as vectors, matrices, their operations and properties, and especially eigenvectors and eigenvalues.

12.2.1 Singular value decomposition and covariance matrices

PCA is based on several mathematical theorems (e.g., Karhunen–Loève theorem) that demonstrate that the space represented by a data matrix Y of size $n \times p$ can be decomposed exactly in rank-1 vectors U and V , such as

$$Y = U S V^T$$

U , S , and V are obtained using a method called singular value decomposition (SVD); S is a diagonal matrix, where all the off-diagonal elements are zeros, and the diagonal elements—called singular values—are ordered such as $\lambda_1 > \lambda_2 > \dots > \lambda_p$.

$$S = \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & \dots & \ddots & \dots \\ 0 & 0 & \dots & \lambda_p \end{bmatrix}$$

U is a $n \times p$ orthonormal matrix: each of the p columns of U represents the eigenvector or principal axes of Y , and they are all orthogonal to each other, i.e., $UU^T = U^T U = I$ (I = identity matrix) and $U^{-1} = U^T$. V is a $p \times p$ orthonormal matrix. Each of the p columns of V represents the eigenvector or principal axes of Y^T , and they are all orthogonal to each other, i.e., $VV^T = V^T V = I$ and $V^{-1} = V^T$.

Fig. 12.2 shows a graphic representation of the decomposition of the matrix Y into a set of rank-1 vectors U and V . The number of vectors is equal to the number of singular values that are nonnull. This number is equal to the rank of matrix Y with $\text{Rank}(Y) \leq \min(n, p)$.

These notations are important because eigenvectors of the data matrix actually represent the principal axes. Each eigenvector has a corresponding eigenvalue. The eigenvalue is a scalar that represents how much variance is explained by the eigenvector when the data are projected on it. The sum of all the eigenvalues corresponds to the total variance of the data.

An alternative way to view the principal components is to use the eigenvectors of the covariance matrices C_p ($C_p = Y^T Y$) and C_n ($C_n = Y Y^T$). The eigenvectors of these two square symmetric matrices are, respectively, the V and U that we obtained before. Singular values S correspond to the square roots of eigenvalues. In fact, we have the following relationships:

$$\begin{aligned} Y &= US V^T \\ Y^T Y V &= V S^2 \rightarrow V = \text{eig}(Y^T Y) = \text{eig}(C_p) \\ Y Y^T U &= U S^2 \rightarrow U = \text{eig}(Y Y^T) = \text{eig}(C_n) \end{aligned}$$

12.2.2 Dimension reduction and variance explained

Once the eigenvectors U and V and the eigenvalue have been calculated using the eigendecomposition, or the SVD approach, we can reconstruct new data Y using fewer components. Below, we reconstruct Y using only the principal components up to the rank r , the r principal components with largest variance explained. Our assumption is that the remaining

components have small contribution to the variance and correspond to noise. The reconstructed Y is equal to

$$Y_{rec} = U(:, 1:r) \text{diag}(S_{[1\dots r]}) V(:, 1:r)^T$$

As PCA is a fully data-driven method with no assumption on the distribution of the noise, there are no straightforward statistical tests to decide on the number of significant components. One way to decide how many components to retain is to compute the cumulative variance explained and decide on a particular threshold, e.g., 80%. The percentage of variance of each component and the cumulative variance explained are equal to

$$\begin{aligned} (\text{percent})Var_i &= \frac{\lambda_i}{\sum \lambda_i} * 100 \\ (\text{cumulative})Var_r &= \sum_{i=1}^r \text{percent } Var_i \end{aligned}$$

12.2.3 Extensions of PCA and other multivariate methods

PCA is a flexible method that can easily be tuned to address different problems and situations. There are different multivariate methods—PLS (partial least squares), MLM (multivariate), CVA (canonical variate analyses)—all of which can be linked to PCA or the SVD methods, as described in [Table 12.1](#) below from [Kherif et al. \(2002\)](#).

PCA can also be applied in the context of regression models ([Fig. 12.3](#)). Regression analysis aims to assess the association between the observations in a data matrix Y and the predictors in a data matrix X . When X and Y are large, using PCA can provide multiple insights into the putative associations. PCA is used in the traditional way on matrix Y ; through decomposition of the covariance matrix, one is able to reveal the principal axes of variation of the data. PCA can also be applied to design matrix X ; this is particularly useful when X contains multiple predictors that are correlated. By performing PCA on the design matrix, we can replace the original design matrix with the reduced uncorrelated principal components. This method is known as principal component regression. A further use of PCA involves the decomposition of the parameters estimates B , which is proportional to the covariance matrix between X and Y . The method is used to find linear combinations of Y that are maximally correlated to linear combinations of X . Finally, the PCA can be used to explore the residuals of the model. If the principal components reveal structured patterns, this could mean that an important variable is missing from the model or that the data contain some artifacts.

TABLE 12.1 Principal component analysis (PCA)—based methods. Most of the multivariate methods can be computed using singular value decomposition (SVD). This table presents methods based on SVD. The second column is the matrix decomposed by SVD. We give also some of the shortcomings associated with each of these methods. In this table, Y are the data, X the linear model, R the orthogonal projector onto the residual space, and Σ the temporal covariance.

Method	Information matrix to decompose by SVD	Pros and cons
PCA	Y	<ul style="list-style-type: none"> • No prior knowledge included • Not data scale invariant
Partial least squares (PLS)	$X^T Y$	<ul style="list-style-type: none"> • Not model and data scale invariant • Problem of temporal correlation
Orthonormalized PLS	$(X^T X)^{-1/2} X^T Y$	<ul style="list-style-type: none"> • Not data scale invariant • Problem of temporal correlation
Canonical variate analysis (CVA)	$(X^T X)^{-1/2} X^T Y (Y^T R Y)^{-1/2}$	<ul style="list-style-type: none"> • Problem of the $(Y^T R Y)^{-1/2}$ computation
SVD-CVA	$Y_k = U_k S_k V_k (X^T X)^{-1/2} X^T Y_k (Y_k^T Y_k)^{-1/2}$	<ul style="list-style-type: none"> • Problem of finding k
Multivariate analysis	$(X^T \Sigma X)^{-1/2} X^T Y$	<ul style="list-style-type: none"> • Not data scale invariant

Adapted from Worsley, K. J., Poline, J. B., Friston, K. J., & Evans, A. C. (1997). Characterizing the response of PET and fMRI data using multivariate linear models. *NeuroImage*, 6(4), 305–319; Kherif, F., Poline, J. B., Flandin, G., Benali, H., Simon, O., Dehaene, S., et al. (2002). Multivariate model specification for fMRI data. *NeuroImage*, 16(4), 1068–1083.

12.2.4 How PCA works using a toy example dataset with dementia patients

12.2.4.1 Data, analysis, and results

This section aims to illustrate how PCA can be applied to neuroimaging data to investigate brain disorders. The data include T1-weighted images acquired from 40 healthy controls (HCs) and 40 Alzheimer's disease (AD) patients using a 3T GE Prisma scanner. Images were preprocessed using Statistical Parametric Mapping software (SPM12). The gray matter volume

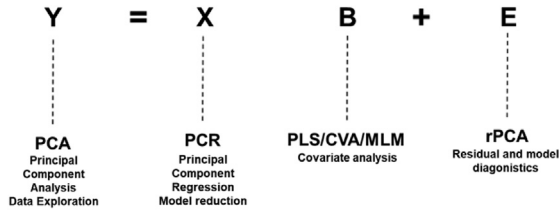


FIGURE 12.3 Principal component analysis (PCA) and model-based approaches: PCA can be applied in the context of a model-based approach such as a linear regression model. PCA can be applied to the raw (mean centered or scaled data) for initial data exploration. PCA can be used to reduce the dimensions of the design matrix X , identifying the principal covariates between X and Y or, finally, identifying potential outliers in the model residuals.

was extracted after segmentation and normalization to the Montreal Neurological Institute (MNI) space. The usual approach for analyzing this dataset is to perform a voxel-wise two sample t -test to identify the voxels where there are significant differences between the two groups (the model includes age, gender, and total intracranial volume as additional confounding variables).

The main steps of PCA are as follows:

1. Load and format the data into a data matrix Y (N voxels \times 80 subjects). Here, each image is presented as a 1D column; thus we have data matrix with 80 columns and N rows.
2. Mean center the data in Y by removing the average.
3. Compute the eigenvectors using the SVD of the matrix $Y = U S V^T$. The columns of U represent the eigenvector of size N voxels, also called eigenimages. Eigenimages can be saved as nifti files, a common format to store neuroimaging data. The columns of V represent the associated eigenvectors of size 80 (equal to the number of subjects). The coefficients of each vector V are called subject loadings and represent the contribution of each subject to the spatial pattern in the corresponding eigenimage.

12.2.4.2 Interpretation

Fig. 12.4 shows the first component or eigenimage on the left, which is very similar to the statistical t -map from the voxel-wise analysis on the right. We expected to observe from the t -map that most of the differences due to AD would be located in the medial temporal lobe (MTL) and in particular in the hippocampus. The first component, which explains most of the variance of the data, not only captures the effects of the disease but also shows more extended brain changes in the other brain regions, such as the cerebellum. We can conclude that, in this dataset, the effects of the disease explained most of the variance above and beyond other effects such as age and gender.

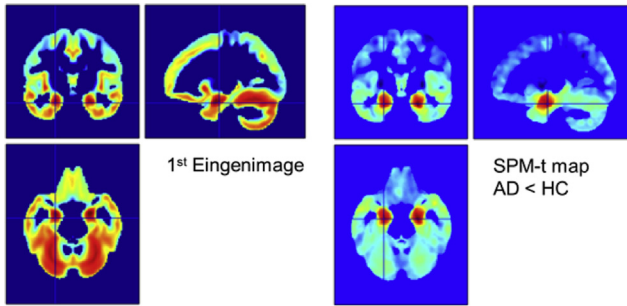


FIGURE 12.4 Pattern of Alzheimer's disease identified by principal component analysis (PCA) or by linear regression. (Left): The first component or eigenimage computed by PCA. (Right): The statistical t-map computed by voxel-wise regression analysis. The crosshairs are centered at the left hippocampus.

To assess how the PCA components discriminate between patients and controls, we can plot the first subject loading against the second subject loading. As Fig. 12.5 shows, in this plot, we can easily discriminate between the patients (coded as 1) and the HCs (coded as 2). The first two subject loadings can be used as features for machine learning tools (KNN, SVM, etc.), the advantage being that vectors have smaller size (1,80) compared to the thousands of voxels in the preprocessed images.

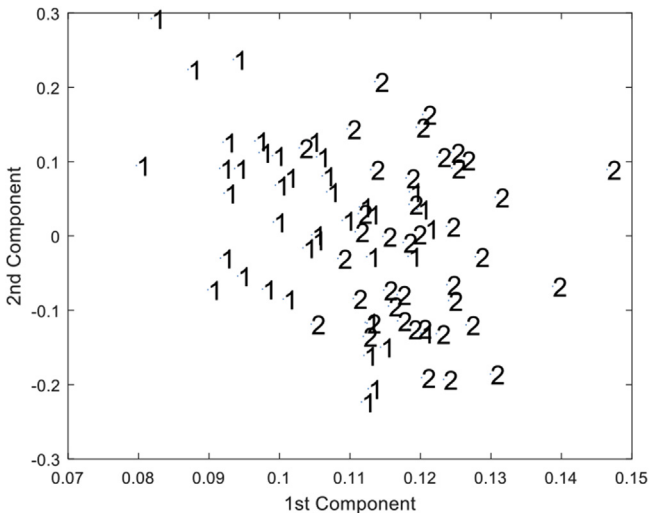


FIGURE 12.5 Principal component analysis subject loadings: the plot shows the first subject loadings against the second subject loadings. The Alzheimer's disease patients are encoded as 1 and the healthy controls are encoded as 2.

12.3 Applications to brain disorders

PCA and its variants (PLS, CVA, MLM) are largely used in studies aiming to identify the neural mechanisms of brain disorders. In particular, an increasing number of studies are using machine learning tools for building classification tools with many using PCA components as input features. In this section, we consider applications of PCA to psychiatric disorders ([Section 12.1.4.1](#)), neurological disorders ([Section 12.1.4.2](#)), and healthy aging ([Section 12.1.4.3](#)).

12.3.1 Psychiatric disorders

Psychiatric disorders are normally diagnosed only on the basis of clinical symptoms. Currently, there are no biomarkers or laboratory tests available to inform the diagnostic and prognostic assessment of psychiatric disorders. An increasing number of researchers are therefore combining neuroimaging data with MLM analyses to develop objective diagnostic and prognostic tools.

In a study by [Kawasaki et al. \(2007\)](#), the authors applied PCA-based MLM linear models to structural magnetic resonance imaging (MRI) data to identify neuroanatomical alternations which could be used to discriminate between patients with schizophrenia and HCs. The data analysis workflow is illustrated in [Fig. 12.6](#), and the results are shown in [Fig. 12.7](#). The results indicated that PCA-based approaches are able to capture the main source of variation associated with schizophrenia and that the pattern captured in the subject loading was predictive of the presence of the disease in an independent dataset. Specifically, when the eigenimage derived from the original cohort was applied to the second cohort, it correctly assigned more than 80% of the healthy subjects and schizophrenia patients.

In a subsequent multisite study, the authors used MLM approach to investigate the neuroanatomical correlates of genetic susceptibility to autism spectrum disorder, schizophrenia, and language and cognitive impairment ([Maillard et al., 2015](#)). The data included MRI imaging data collected via 3T scanners and whole genome arrays confirming either a recurrent deletion or duplication of the genetic 16p11.2 breakpoint 4 to 5 copy number variants from people with deletion ($n = 14$), duplication ($n = 17$), and controls ($n = 23$). Specifically, the authors used the PCA-based MLM to quantify the effect of either a recurrent deletion or duplication of the BP4-BP5 region on brain structure. Using the voxel-wise univariate method to test for a dosage-dependent gene effect, the authors found significant differences in the insula (i.e., deletion > control > duplication). A number of additional regions were affected by either the deletion or the duplication; these included the calcarine cortex and transverse temporal gyrus (deletion > control), the

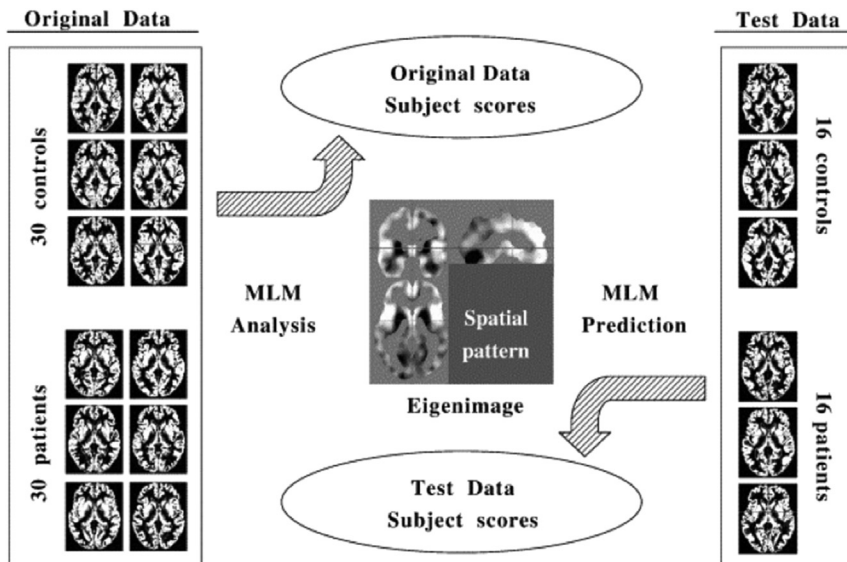


FIGURE 12.6 Principal component analysis (PCA) for discriminative analyses. An eigenimage and the subject loading are calculated using the PCA-based multivariate (MLM) analysis within the original data. The resulting first eigenimage is then used as a predictor for classification of separate test data using the subject loadings of the new data. Taken from Kawasaki, Y., Suzuki, M., Kherif, F., Takahashi, T., Zhou, S. Y., Nakamura, K., et al. (2007). Multivariate voxel-based morphometry successfully differentiates schizophrenia patients from healthy controls. *NeuroImage*, 34(1), 235–242. Epub 2006 Oct 11. PubMed PMID: 17045492.

superior and middle temporal gyri (deletion < control), and the caudate and hippocampus (control > duplication). In a subsequent analysis, PCA-based MLM was used to decompose the cross-correlation matrix between gene expression and regional brain volume. This approach involves using all mRNA levels and all voxels to determine the best representations within each type of data that explain maximum covariance between gene expression and regional brain volume across all individuals. The eigenimage, represented in the brain space, showed the degree of contribution (either positive or negative) of each voxel to the correlation mapping. Similarly, the gene loadings (positive or negative) show the contribution of each gene to the correlation mapping.

12.3.2 Neurological disorders

Neurological diseases such as Parkinson's disease (PD) and AD affect a large proportion of the aging population (more than 50% after the age of 65 years). In the recent years, there has been a surge in studies using complex MLM analyses and machine learning to identify the biological signatures of these diseases and develop diagnostic and prognostic models.

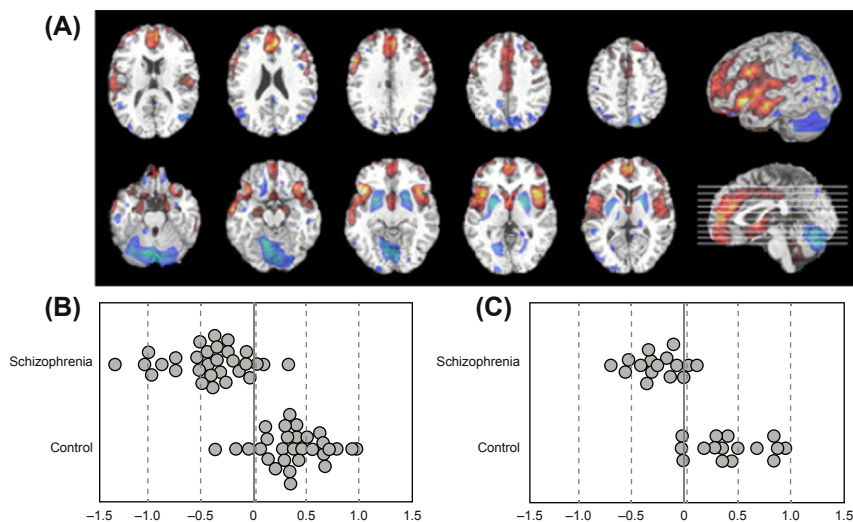


FIGURE 12.7 Principal component analysis (PCA)-based MLM results. (A) Eigenimage with the 30% largest loadings for positive (in red) (dark gray in print version) and negative (in blue) (gray in print version) loadings. (B) Subject loadings in the training. (C) Subjects loadings in the test set. Adapted from Kawasaki, Y., Suzuki, M., Kherif, F., Takahashi, T., Zhou, S. Y., Nakamura, K., et al. (2007). Multivariate voxel-based morphometry successfully differentiates schizophrenia patients from healthy controls. *NeuroImage*, 34(1), 235–242. Epub 2006 Oct 11. PubMed PMID: 17045492.

Xiao et al. (2014), for example, used PCA to explain individual differences in the position and size of the basal ganglia substructures in PD patients. For all structures of interest, the first five components accounted for 95%–98% of the total variability, whereas the first 10 components accounted for roughly 99% of the total variability; within the basal ganglia, the red nucleus was less variable than the substantia nigra and the subthalamic nucleus. These findings are of great interest because the vast majority of neuroimaging studies use standard atlas coordinates to target the basal ganglia, therefore neglecting the large amount of individual variability. In addition, post hoc correlation analyses showed that principal components were associated with the current disease manifestations and other markers of disease progression.

Most of the studies of neurological disorders use PCA as a preprocessing approach for dimensionality reduction. We saw previously that an important choice is the number of components to retain. This question was also addressed in an interesting study (Markiewicz, Matthews, Declerck, Herholz, & Alzheimer's Disease Neuroimaging Initiative (ADNI), 2011) which compared patients with AD and HCs using a large-scale PET data from the Alzheimer's Disease Neuroimaging Initiative (ADNI) cohort (<http://adni.loni.usc.edu/>). The results—replicated in two independent datasets—suggested that the first four components captured most of the variance in the data.

PCA-based methods are also useful for investigating the interaction between multiple high-dimensional modalities. [Zufferey et al. \(2017\)](#) used hierarchical MLM linear models to quantify the interaction between personality traits, state of cognitive impairment, and MRI measures (gray matter brain volume and gray matter mean water diffusion) in the MTL. The MLM linear model showed the interaction between cognitive state and personality traits predicting MTL abnormalities was mainly driven by neuroticism and its facets (anxiety, depression, and stress) and was associated with right–left asymmetry and an anterior-to-posterior gradient in the MTL ([Fig. 12.8](#)).

12.3.3 Healthy aging and disease process

To conclude this section illustrating possible applications of PCA to brain disorders, we discuss what is probably the most important factor in disease modeling—age. Age is the main risk factor for the most neurological diseases such as AD and PD. This means it can be challenging to disentangle the effects of a disease of interest from the effects of age. In most studies, for example, age is treated as confounding factor which involves removing its effect from the dependent variable. The problem with this approach is that age tends to correlate with clinical variables of interest (e.g., clinical and cognitive symptoms, duration of illness, duration of medication) and therefore modeling age as a confounding factor comes with a high risk of removing disease-related variability from the data. To address this problem, it is important to develop robust models of “healthy brain aging,” which can be used to accurately discriminate between natural changes and disease-related changes. PCA-based methods provide a useful set of tools for elucidating the effect of aging on the brain and how this varies according to the presence or absence of disease. [Franke, Ziegler, Klöppel, Gaser, and Alzheimer’s Disease Neuroimaging Initiative \(2010\)](#) combined PCA with relevance vector machine to build a model that predicts chronological age based on T1-weighted scans. The authors suggested using the difference between the actual age and the predicted brain age as a biomarker for brain health. While this investigation used a single type of data (i.e., T1-weighted scans), a growing number of studies are employing multimodal data to investigate brain age. The so-called biological age score can be computed from multimodal data (cognitive, physical, physiological, or biochemical data) using principal component approach as detailed in [Nakamura and Miyao, \(2007\)](#). A similar multimodal approach was followed by [Draganski et al. \(2011\)](#) who computed a MLM correlate of biological age using multiparametric maps from quantitative MRI contrasts ([Fig. 12.9](#)).

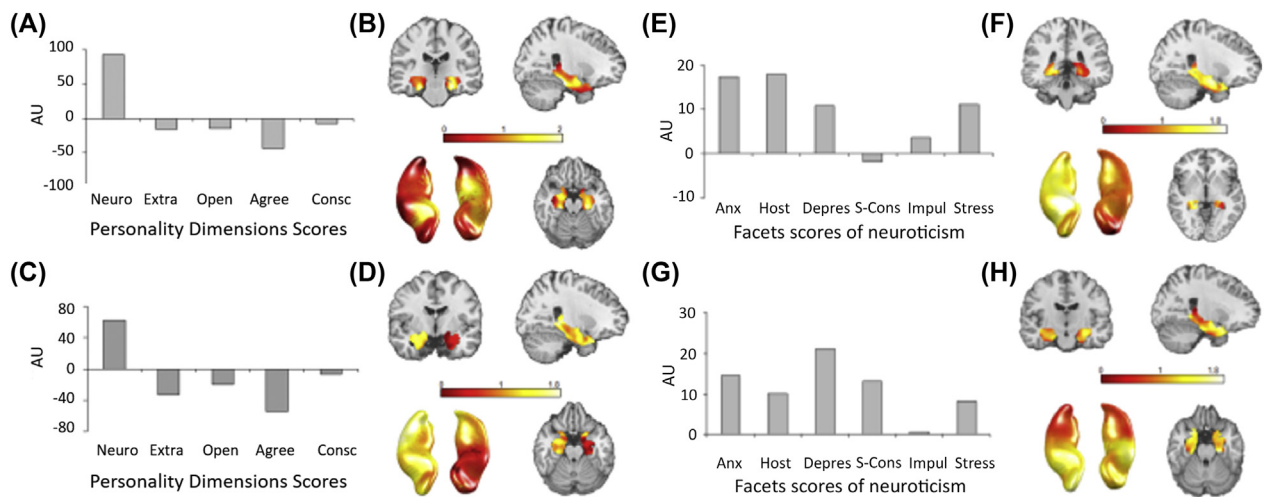


FIGURE 12.8 Principal component analysis (PCA)-based MLM analysis of personality profile at domain level: (A) First eigen component ($P < .05$) and (B) the associated spatial distribution within the search volume of interest for gray matter volume; (C) First Eigen-component ($P < .05$); (D) the associated spatial distribution within the search volume of interest for gray matter mean diffusivity and (E,F,G,H) show the anatomical differences related to personality facets of neuroticism and cognitive state. *Neuro*, neuroticism; *Extra*, extraversion; *Open*, openness; *Agree*, agreeableness; *Consc*, conscientiousness. Taken from Zufferey, V., Donati, A., Popp, J., Meuli, R., Rossier, J., Frackowiak, R., et al. (2017). Neuroticism, depression, and anxiety traits exacerbate the state of cognitive impairment and hippocampal vulnerability to Alzheimer's disease. *Alzheimer's and Dementia: Diagnosis, Assessment and Disease Monitoring*, 7, 107–114.

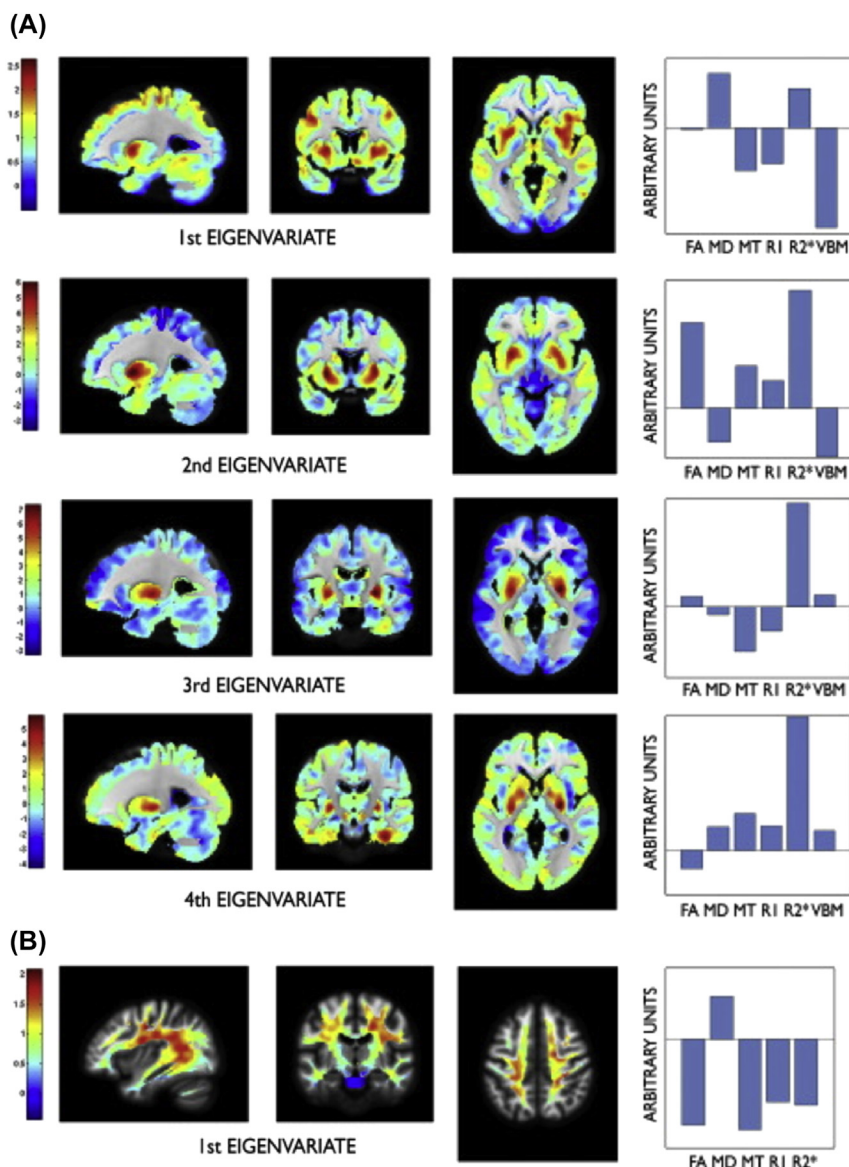


FIGURE 12.9 Spatial patterns of statistically significant eigenvars ($P < .05$, Lawley–Hotelling trace) related to age in gray (A) and white matter (B) using multivariate linear model of age. The first and second eigenvars in the gray matter explained 80% of the variance, whereas the first eigenvariate in the white matter explained 86% of the variance. The specific pattern of directionality and magnitude of change in the parameter maps is presented in the bar graphs on the right-hand side. FA, fractional anisotropy; MD, mean diffusivity; MT, magnetization transfer, R1, R2*; VBM, voxel-based morphometry (i.e., gray matter volume). Taken from Draganski, B., Ashburner, J., Hutton, C., Kherif, F., Frackowiak, R. S., Helms, G., et al. (2011). Regional specificity of MRI contrast parameter changes in normal ageing revealed by voxel-based quantification (VBQ). *NeuroImage*, 55(4), 1423–1434.

12.4 Conclusion

In this chapter, we have presented the method of PCA and discussed some of its potential applications to brain disorders. PCA aims to find a compact representation of the data by projecting them into new orthogonal axes. In contrast with most of the machine learning models discussed in the present book, it is a very old method, first proposed by Karl Pearson in 1901 (Pearson, 1901). More than 100 years later, the method is still used across several domains including, among others, physics, biology, and psychology.

It is important to acknowledge that PCA also suffers from a number of limitations: the method is linear, and it can be difficult to determine the statistical significance of the components. In addition, sometimes the results can be challenging to interpret, for example, when the orthogonality constraints do not lead to biologically plausible components. On the other hand, it has the advances in storage, memory, distributed computing, and novel eigendecomposition algorithms, PCA can be scaled to ever increasing datasets. In addition, with the increasing popularity of machine learning approaches, PCA is being widely used to generate a reduced set of features that can be used instead of the original data. In light of its scalability and ability to simplify large and complex data, therefore, we can expect PCA to continue to be a popular method in the era of Big Data.

12.5 Key points

- PCA is a powerful method for capturing the variance in the data including, but not limited to, individual differences.
- The main aim of PCA is the identification of a reduced set of features that represent the original data in a lower-dimensional subspace with a minimal loss of information.
- PCA can be tuned depending on the question of interest, from data exploration to model prediction and to model diagnostic.
- PCA has multiple potential applications in the investigation of psychiatric and neurological disorders.
- PCA has a number of limitations—for example, the method is linear, it can be difficult to determine statistical significance, and the results can be challenging to interpret.

References

- Draganski, B., Ashburner, J., Hutton, C., Kherif, F., Frackowiak, R. S., Helms, G., et al. (2011). Regional specificity of MRI contrast parameter changes in normal ageing revealed by voxel-based quantification (VBQ). *NeuroImage*, 55(4), 1423–1434.

- Franke, K., Ziegler, G., Klöppel, S., Gaser, C., & Alzheimer's Disease Neuroimaging Initiative. (2010). Estimating the age of healthy subjects from T1-weighted MRI scans using kernel methods: Exploring the influence of various parameters. *NeuroImage*, 50(3), 883–892.
- Jolliffe, I. T. (1986). *Principal component analysis*. 1986. New York: Springer-verlag.
- Kawasaki, Y., Suzuki, M., Kherif, F., Takahashi, T., Zhou, S. Y., Nakamura, K., et al. (2007). Multivariate voxel-based morphometry successfully differentiates schizophrenia patients from healthy controls. *NeuroImage*, 34(1), 235–242. Epub 2006 Oct 11. PubMed PMID: 17045492.
- Kherif, F., Poline, J. B., Flandin, G., Benali, H., Simon, O., Dehaene, S., et al. (2002). Multivariate model specification for fMRI data. *NeuroImage*, 16(4), 1068–1083.
- Maillard, A. M., Ruef, A., Pizzagalli, F., Migliavacca, E., Hippolyte, L., Adaszewski, S., et al. (2015). The 16p11. 2 locus modulates brain structures common to autism, schizophrenia and obesity. *Molecular Psychiatry*, 20(1), 140.
- Markiewicz, P. J., Matthews, J. C., Declerck, J., Herholz, K., & Alzheimer's Disease Neuroimaging Initiative (ADNI). (2011). Verification of predicted robustness and accuracy of multivariate analysis. *NeuroImage*, 56(3), 1382–1385.
- Nakamura, E., & Miyao, K. (2007). A method for identifying biomarkers of aging and constructing an index of biological age in humans. *The Journals of Gerontology Series A: Biological Sciences and Medical Sciences*, 62(10), 1096–1105.
- Pearson, K. (1901). LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11), 559–572.
- Worsley, K. J., Poline, J. B., Friston, K. J., & Evans, A. C. (1997). Characterizing the response of PET and fMRI data using multivariate linear models. *NeuroImage*, 6(4), 305–319.
- Xiao, Y., Jannin, P., D'Albis, T., Guizard, N., Haegelen, C., Lalys, F., et al. (2014). Investigation of morphometric variability of subthalamic nucleus, red nucleus, and substantia nigra in advanced Parkinson's disease patients using automatic segmentation and PCA-based analysis. *Human Brain Mapping*, 35(9), 4330–4344.
- Zufferey, V., Donati, A., Popp, J., Meuli, R., Rossier, J., Frackowiak, R., et al. (2017). Neuroticism, depression, and anxiety traits exacerbate the state of cognitive impairment and hippocampal vulnerability to Alzheimer's disease. *Alzheimer's and Dementia: Diagnosis, Assessment and Disease Monitoring*, 7, 107–114.