

STATLAB WORKSHOP

PRINCIPAL COMPONENTS ANALYSIS WITH R

W. D. BRINDA

YALE UNIVERSITY

01/28/2019



- 1 Theory
 - Dimension reduction
 - The PCA algorithm

- 1 Theory
 - Dimension reduction
 - The PCA algorithm

- 2 Implementation in R
 - Running example code together
 - Exercises to try for yourself

THEORY

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

There are a variety of reasons why one might want to do this. For example,

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

There are a variety of reasons why one might want to do this. For example,

- Visualization, e.g. reduce many quantitative variables down to a few, then draw scatterplots

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

There are a variety of reasons why one might want to do this. For example,

- Visualization, e.g. reduce many quantitative variables down to a few, then draw scatterplots
- Meaningful interpretations of relationships among variables

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

There are a variety of reasons why one might want to do this. For example,

- Visualization, e.g. reduce many quantitative variables down to a few, then draw scatterplots
- Meaningful interpretations of relationships among variables
- Decreasing variance of estimation (while likely increasing bias)

DIMENSION REDUCTION

Dimension reduction means applying some procedure to a dataset that results in a new dataset with fewer columns. The new dataset can be considered a lower-dimensional "approximation" of the original one.

There are a variety of reasons why one might want to do this. For example,

- Visualization, e.g. reduce many quantitative variables down to a few, then draw scatterplots
- Meaningful interpretations of relationships among variables
- Decreasing variance of estimation (while likely increasing bias)
- A technique that you want to use is computationally impossible or intractable in the data's original dimension

BIAS-VARIANCE TRADE-OFF

Suppose $\hat{\theta}$ is an estimator for an unknown parameter vector $\theta \in \mathbb{R}^d$. The expected squared error of $\hat{\theta}$ is equal to its "squared bias" plus its "variance":

BIAS-VARIANCE TRADE-OFF

Suppose $\hat{\theta}$ is an estimator for an unknown parameter vector $\theta \in \mathbb{R}^d$. The expected squared error of $\hat{\theta}$ is equal to its "squared bias" plus its "variance":

$$\mathbb{E}\|\hat{\theta} - \theta\|^2 = \|\mathbb{E}\hat{\theta} - \theta\|^2 + \mathbb{E}\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2.$$

BIAS-VARIANCE TRADE-OFF

Suppose $\hat{\theta}$ is an estimator for an unknown parameter vector $\theta \in \mathbb{R}^d$. The expected squared error of $\hat{\theta}$ is equal to its "squared bias" plus its "variance":

$$\mathbb{E}\|\hat{\theta} - \theta\|^2 = \|\mathbb{E}\hat{\theta} - \theta\|^2 + \mathbb{E}\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2.$$

(A trade-off between bias and variance holds not only for estimation of an unknown parameter or distribution, but also for prediction of new data.)

BIAS-VARIANCE TRADE-OFF

Suppose $\hat{\theta}$ is an estimator for an unknown parameter vector $\theta \in \mathbb{R}^d$. The expected squared error of $\hat{\theta}$ is equal to its "squared bias" plus its "variance":

$$\mathbb{E}\|\hat{\theta} - \theta\|^2 = \|\mathbb{E}\hat{\theta} - \theta\|^2 + \mathbb{E}\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2.$$

(A trade-off between bias and variance holds not only for estimation of an unknown parameter or distribution, but also for prediction of new data.)

In general, "simpler" models result in larger squared bias but smaller variance. An essential part of the practice of statistics and machine learning is choosing a procedure that balances bias and variance such that their sum is nearly as small as possible.

BIAS-VARIANCE TRADE-OFF

Suppose $\hat{\theta}$ is an estimator for an unknown parameter vector $\theta \in \mathbb{R}^d$. The expected squared error of $\hat{\theta}$ is equal to its "squared bias" plus its "variance":

$$\mathbb{E}\|\hat{\theta} - \theta\|^2 = \|\mathbb{E}\hat{\theta} - \theta\|^2 + \mathbb{E}\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2.$$

(A trade-off between bias and variance holds not only for estimation of an unknown parameter or distribution, but also for prediction of new data.)

In general, "simpler" models result in larger squared bias but smaller variance. An essential part of the practice of statistics and machine learning is choosing a procedure that balances bias and variance such that their sum is nearly as small as possible.

Dimension reduction techniques tend to increase squared bias and decrease variance.

There are two common *linear* transformations for reducing the dimension of a set of quantitative variables.

There are two common *linear* transformations for reducing the dimension of a set of quantitative variables.

- Principal components analysis (PCA) "spreads out the *points* as much as possible"

There are two common *linear* transformations for reducing the dimension of a set of quantitative variables.

- Principal components analysis (PCA) "spreads out the *points* as much as possible"
- Linear discriminant analysis (LDA) "spreads out the *groups* as much as possible" (requires a categorical variable)

There are two common *linear* transformations for reducing the dimension of a set of quantitative variables.

- Principal components analysis (PCA) "spreads out the *points* as much as possible"
- Linear discriminant analysis (LDA) "spreads out the *groups* as much as possible" (requires a categorical variable)

PCA is our focus, but we'll also see LDA in one of the R code examples later.

A "SIMPLE" EXPLANATION

In the following slides, I've tried to present a thorough explanation of the PCA algorithm *as simply as possible...* that doesn't mean that it's *simple*.

A "SIMPLE" EXPLANATION

In the following slides, I've tried to present a thorough explanation of the PCA algorithm *as simply as possible...* that doesn't mean that it's *simple*.

To be frank, it's probably more than you can process in a single sitting. Just do your best to follow along and see how much you can make sense of. Then, if you'd like to really understand the details, take your time looking over these notes again later.

SPREADING OUT THE DATA POINTS

Consider a data frame with n observations of d quantitative variables. One could try to imagine the high-dimensional scatterplot of the n observations as data points in \mathbb{R}^d . For simplicity, just think of an ordinary scatterplot of two variables.

SPREADING OUT THE DATA POINTS

Consider a data frame with n observations of d quantitative variables. One could try to imagine the high-dimensional scatterplot of the n observations as data points in \mathbb{R}^d . For simplicity, just think of an ordinary scatterplot of two variables.

What exactly does it mean to "spread out" the points as much as possible? Can you think of a common statistic for quantifying how spread out a set of numbers is?

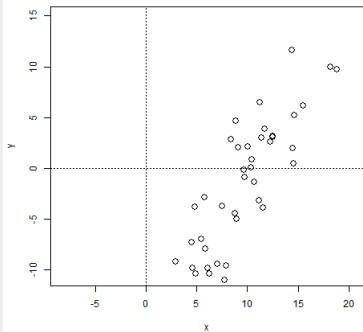
SPREADING OUT THE DATA POINTS

Consider a data frame with n observations of d quantitative variables. One could try to imagine the high-dimensional scatterplot of the n observations as data points in \mathbb{R}^d . For simplicity, just think of an ordinary scatterplot of two variables.

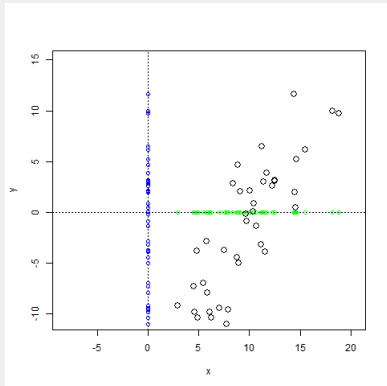
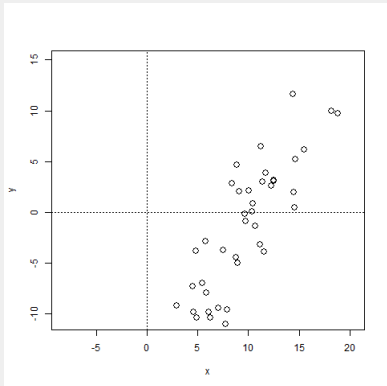
What exactly does it mean to "spread out" the points as much as possible? Can you think of a common statistic for quantifying how spread out a set of numbers is?

How about *variance*, the average squared difference from the mean. Each variable has a variance, and some variables will have larger variances than others.

A CLOUD OF DATA POINT AND THEIR COORDINATES



A CLOUD OF DATA POINT AND THEIR COORDINATES



ALTERNATIVE BASES

Consider the first observation $\mathbf{x}_1 = (x_{1,1}, \dots, x_{1,d})$. It's the sum of its orthogonal projections onto the coordinate axes:

$$\mathbf{x}_1 = x_{1,1}\mathbf{e}_1 + \dots + x_{1,d}\mathbf{e}_d$$

where \mathbf{e}_j is the j th standard basis vector (it has a 1 as its j th entry and zeros elsewhere).

ALTERNATIVE BASES

Consider the first observation $\mathbf{x}_1 = (x_{1,1}, \dots, x_{1,d})$. It's the sum of its orthogonal projections onto the coordinate axes:

$$\mathbf{x}_1 = x_{1,1}\mathbf{e}_1 + \dots + x_{1,d}\mathbf{e}_d$$

where \mathbf{e}_j is the j th standard basis vector (it has a 1 as its j th entry and zeros elsewhere).

However, \mathbf{x}_1 is also the sum of its orthogonal projections onto *any* orthonormal basis $\mathbf{u}_1, \dots, \mathbf{u}_d$. The general formula is simply

$$\mathbf{x}_1 = (\mathbf{x}_1^T \mathbf{u}_1)\mathbf{u}_1 + \dots + (\mathbf{x}_1^T \mathbf{u}_d)\mathbf{u}_d.$$

ALTERNATIVE BASES

Consider the first observation $\mathbf{x}_1 = (x_{1,1}, \dots, x_{1,d})$. It's the sum of its orthogonal projections onto the coordinate axes:

$$\mathbf{x}_1 = x_{1,1}\mathbf{e}_1 + \dots + x_{1,d}\mathbf{e}_d$$

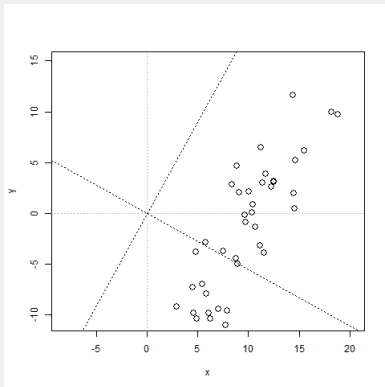
where \mathbf{e}_j is the j th standard basis vector (it has a 1 as its j th entry and zeros elsewhere).

However, \mathbf{x}_1 is also the sum of its orthogonal projections onto *any* orthonormal basis $\mathbf{u}_1, \dots, \mathbf{u}_d$. The general formula is simply

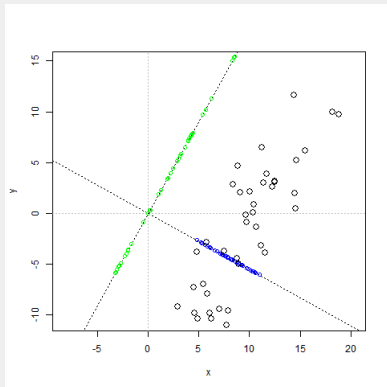
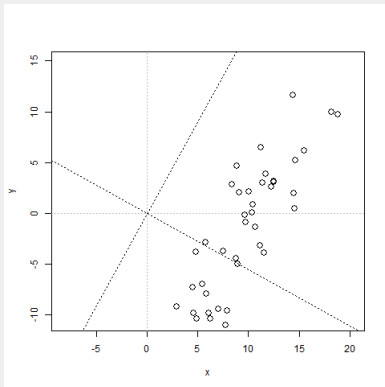
$$\mathbf{x}_1 = (\mathbf{x}_1^T \mathbf{u}_1)\mathbf{u}_1 + \dots + (\mathbf{x}_1^T \mathbf{u}_d)\mathbf{u}_d.$$

The other observations can be represented likewise. Thus, the *coordinates* of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to a basis vector \mathbf{u} are $\mathbf{x}_1^T \mathbf{u}, \dots, \mathbf{x}_n^T \mathbf{u}$.

ALTERNATIVE AXES AND COORDINATES



ALTERNATIVE AXES AND COORDINATES



PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

- The average of the coordinates with respect to any basis vector \mathbf{u} is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ where \mathbb{X} is the matrix whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$.

PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

- The average of the coordinates with respect to any basis vector \mathbf{u} is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ where \mathbb{X} is the matrix whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .

PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

- The average of the coordinates with respect to any basis vector \mathbf{u} is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ where \mathbb{X} is the matrix whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

- The average of the coordinates with respect to any basis vector \mathbf{u} is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ where \mathbb{X} is the matrix whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ (with respect to any basis vector) is exactly the same as the variance of the coordinates of the *centered* vectors $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$.

PREVIEW OF PCA ARGUMENT

Assuming the mean of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the zero vector:

- The average of the coordinates with respect to any basis vector \mathbf{u} is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ where \mathbb{X} is the matrix whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ (with respect to any basis vector) is exactly the same as the variance of the coordinates of the *centered* vectors $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$.

Let $\tilde{\mathbb{X}}$ be the matrix whose rows are $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$, and let $\frac{1}{n} \tilde{\mathbb{X}}^T \tilde{\mathbb{X}}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n$. The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_1 . Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

AVERAGE OF COORDINATES

Suppose that the mean $\frac{1}{n} \sum_i \mathbf{x}_i$ is equal to the zero vector $\mathbf{0}$, that is, the vector with zero as every entry.

AVERAGE OF COORDINATES

Suppose that the mean $\frac{1}{n} \sum_i \mathbf{x}_i$ is equal to the zero vector $\mathbf{0}$, that is, the vector with zero as every entry.

Then the average of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to a unit vector \mathbf{u} is

$$\begin{aligned} \frac{1}{n}(\mathbf{x}_1^T \mathbf{u} + \dots + \mathbf{x}_n^T \mathbf{u}) &= \frac{1}{n} \underbrace{(\mathbf{x}_1 + \dots + \mathbf{x}_n)^T}_{\mathbf{0}} \mathbf{u} \\ &= 0 \end{aligned}$$

VARIANCE OF COORDINATES

If $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , then the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to a unit vector \mathbf{u} is simply the average of the squared coordinates

$$\frac{1}{n}[(\mathbf{x}_1^T \mathbf{u})^2 + \dots + (\mathbf{x}_n^T \mathbf{u})^2].$$

VARIANCE OF COORDINATES

If $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , then the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to a unit vector \mathbf{u} is simply the average of the squared coordinates

$$\frac{1}{n}[(\mathbf{x}_1^T \mathbf{u})^2 + \dots + (\mathbf{x}_n^T \mathbf{u})^2].$$

Let's derive a more useful representation of this quantity in terms of the matrix \mathbb{X} whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$. The coordinates $\mathbf{x}_1^T \mathbf{u}, \dots, \mathbf{x}_n^T \mathbf{u}$ are the entries of the vector

$$\mathbb{X}\mathbf{u} = \begin{bmatrix} - & \mathbf{x}_1 & - \\ & \vdots & \\ - & \mathbf{x}_n & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{u} \\ | \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \mathbf{u} \\ \vdots \\ \mathbf{x}_n^T \mathbf{u} \end{bmatrix}.$$

VARIANCE OF COORDINATES

If $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , then the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to a unit vector \mathbf{u} is simply the average of the squared coordinates

$$\frac{1}{n} [(\mathbf{x}_1^T \mathbf{u})^2 + \dots + (\mathbf{x}_n^T \mathbf{u})^2].$$

Let's derive a more useful representation of this quantity in terms of the matrix \mathbb{X} whose rows are $\mathbf{x}_1, \dots, \mathbf{x}_n$. The coordinates $\mathbf{x}_1^T \mathbf{u}, \dots, \mathbf{x}_n^T \mathbf{u}$ are the entries of the vector

$$\mathbb{X}\mathbf{u} = \begin{bmatrix} - & \mathbf{x}_1 & - \\ & \vdots & \\ - & \mathbf{x}_n & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{u} \\ | \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \mathbf{u} \\ \vdots \\ \mathbf{x}_n^T \mathbf{u} \end{bmatrix}.$$

Therefore, the variance of the coordinates can be represented

$$\frac{1}{n} \sum_i (\mathbf{x}_i^T \mathbf{u})^2 = \frac{1}{n} \|\mathbb{X}\mathbf{u}\|^2 = \frac{1}{n} (\mathbb{X}\mathbf{u})^T (\mathbb{X}\mathbf{u}) = \mathbf{u}^T \left(\frac{1}{n} \mathbb{X}^T \mathbb{X} \right) \mathbf{u}.$$

MAXIMIZING VARIANCE

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , we've realized that the variance of their coordinates with respect to any unit vector \mathbf{u} equals the *quadratic form* $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.

MAXIMIZING VARIANCE

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , we've realized that the variance of their coordinates with respect to any unit vector \mathbf{u} equals the *quadratic form* $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.

It turns out that this is easy to maximize via a *spectral decomposition*. The symmetric matrix $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ can be rewritten in terms of its eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ and orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$:

$$\frac{1}{n} \mathbb{X}^T \mathbb{X} = \lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \dots + \lambda_d \mathbf{q}_d \mathbf{q}_d^T.$$

MAXIMIZING VARIANCE

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , we've realized that the variance of their coordinates with respect to any unit vector \mathbf{u} equals the *quadratic form* $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.

It turns out that this is easy to maximize via a *spectral decomposition*. The symmetric matrix $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ can be rewritten in terms of its eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ and orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$:

$$\frac{1}{n} \mathbb{X}^T \mathbb{X} = \lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \dots + \lambda_d \mathbf{q}_d \mathbf{q}_d^T.$$

Therefore

$$\begin{aligned} \mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u} &= \mathbf{u}^T (\lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \dots + \lambda_d \mathbf{q}_d \mathbf{q}_d^T) \mathbf{u} \\ &= \lambda_1 (\mathbf{u}^T \mathbf{q}_1)^2 + \dots + \lambda_d (\mathbf{u}^T \mathbf{q}_d)^2. \end{aligned}$$

The numbers $\mathbf{u}^T \mathbf{q}_1, \dots, \mathbf{u}^T \mathbf{q}_d$ are exactly the coordinates of \mathbf{u} with respect to the orthonormal basis $\mathbf{q}_1, \dots, \mathbf{q}_d$. Because \mathbf{u} is a unit vector, their squares must sum to 1 (Parseval's identity). Therefore, $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is a weighted average of the eigenvalues $\lambda_1, \dots, \lambda_d$. Its maximum possible value of λ_1 is achieved if \mathbf{u} equals \mathbf{q}_1 .

MAXIMIZING VARIANCE

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean \mathbf{o} , we've realized that the variance of their coordinates with respect to any unit vector \mathbf{u} equals the *quadratic form* $\mathbf{u}^T(\frac{1}{n}\mathbb{X}^T\mathbb{X})\mathbf{u}$.

It turns out that this is easy to maximize via a *spectral decomposition*. The symmetric matrix $\frac{1}{n}\mathbb{X}^T\mathbb{X}$ can be rewritten in terms of its eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ and orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$:

$$\frac{1}{n}\mathbb{X}^T\mathbb{X} = \lambda_1\mathbf{q}_1\mathbf{q}_1^T + \dots + \lambda_d\mathbf{q}_d\mathbf{q}_d^T.$$

Therefore

$$\begin{aligned}\mathbf{u}^T(\frac{1}{n}\mathbb{X}^T\mathbb{X})\mathbf{u} &= \mathbf{u}^T(\lambda_1\mathbf{q}_1\mathbf{q}_1^T + \dots + \lambda_d\mathbf{q}_d\mathbf{q}_d^T)\mathbf{u} \\ &= \lambda_1(\mathbf{u}^T\mathbf{q}_1)^2 + \dots + \lambda_d(\mathbf{u}^T\mathbf{q}_d)^2.\end{aligned}$$

The numbers $\mathbf{u}^T\mathbf{q}_1, \dots, \mathbf{u}^T\mathbf{q}_d$ are exactly the coordinates of \mathbf{u} with respect to the orthonormal basis $\mathbf{q}_1, \dots, \mathbf{q}_d$. Because \mathbf{u} is a unit vector, their squares must sum to 1 (Parseval's identity). Therefore, $\mathbf{u}^T(\frac{1}{n}\mathbb{X}^T\mathbb{X})\mathbf{u}$ is a weighted average of the eigenvalues $\lambda_1, \dots, \lambda_d$. Its maximum possible value of λ_1 is achieved if \mathbf{u} equals \mathbf{q}_1 .

By similar reasoning, among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

CENTERING

Define the *mean vector* $\bar{\mathbf{x}}$ to be the vector whose first entry is the average of the first entries of $\mathbf{x}_1, \dots, \mathbf{x}_n$, and so on. The *centered* points $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$ have two key properties.

Define the *mean vector* $\bar{\mathbf{x}}$ to be the vector whose first entry is the average of the first entries of $\mathbf{x}_1, \dots, \mathbf{x}_n$, and so on. The *centered* points $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$ have two key properties.

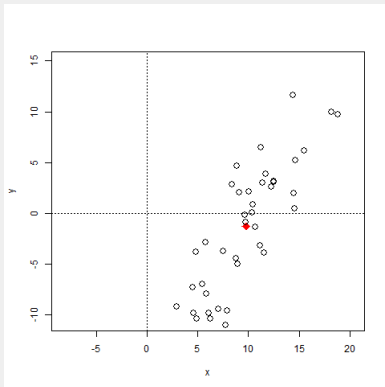
- The coordinates of the centered points with respect to any basis vector have the same variance as the coordinates of the original points with respect to that basis vector.
(Variances don't change when you translate the whole cloud of points by same vector.)

CENTERING

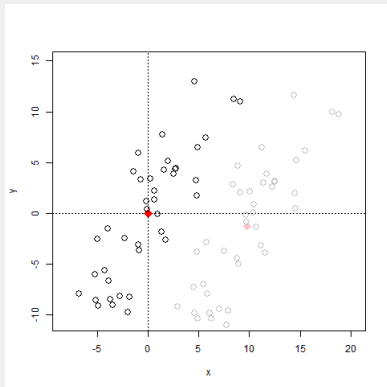
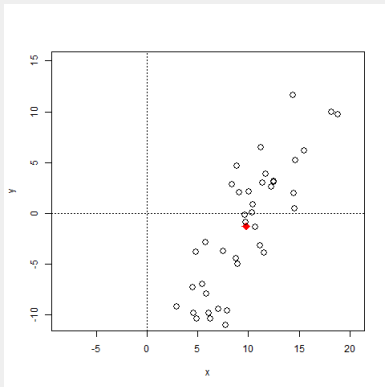
Define the *mean vector* $\bar{\mathbf{x}}$ to be the vector whose first entry is the average of the first entries of $\mathbf{x}_1, \dots, \mathbf{x}_n$, and so on. The *centered* points $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$ have two key properties.

- The coordinates of the centered points with respect to any basis vector have the same variance as the coordinates of the original points with respect to that basis vector. (Variances don't change when you translate the whole cloud of points by same vector.)
- The centered points have mean \mathbf{o} .
$$\left(\frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}})\right) = \left(\frac{1}{n} \sum_i \mathbf{x}_i\right) - \bar{\mathbf{x}} = \mathbf{o}.$$

CENTERING THE DATA POINTS



CENTERING THE DATA POINTS



The variance-maximizing basis vector for $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the same as the variance-maximizing basis vector for the centered points, which we know how to find since they have mean $\mathbf{0}$.

The variance-maximizing basis vector for $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the same as the variance-maximizing basis vector for the centered points, which we know how to find since they have mean $\mathbf{0}$.

Let $\tilde{\mathbf{X}}$ be the matrix whose rows are $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$, and let $\frac{1}{n} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_1 . Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

The variance-maximizing basis vector for $\mathbf{x}_1, \dots, \mathbf{x}_n$ is the same as the variance-maximizing basis vector for the centered points, which we know how to find since they have mean $\mathbf{0}$.

Let $\tilde{\mathbf{X}}$ be the matrix whose rows are $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$, and let $\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_1 . Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

Notice that $\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ is the covariance matrix of the data points.

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

- The average of the coordinates is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

- The average of the coordinates is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

- The average of the coordinates is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

- The average of the coordinates is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ (with respect to any basis vector) is exactly the same as the variance of the coordinates of the *centered* vectors $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$.

RECAP OF PCA ARGUMENT

Assuming $\mathbf{x}_1, \dots, \mathbf{x}_n$ are centered:

- The average of the coordinates is zero, so their variance is equal to the average of the squared coordinates which can be represented $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$.
- Let $\frac{1}{n} \mathbb{X}^T \mathbb{X}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. Spectral decomposition shows that the quadratic form $\mathbf{u}^T (\frac{1}{n} \mathbb{X}^T \mathbb{X}) \mathbf{u}$ is maximized if \mathbf{u} equals \mathbf{q}_1 .
- Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ (with respect to any basis vector) is exactly the same as the variance of the coordinates of the *centered* vectors $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$.

Let $\tilde{\mathbb{X}}$ be the matrix whose rows are $\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}$, and let $\frac{1}{n} \tilde{\mathbb{X}}^T \tilde{\mathbb{X}}$ have eigenvalues $\lambda_1 \geq \dots \geq \lambda_d$ with orthonormal eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_d$. The variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_1 . Among all unit vectors that are orthogonal to \mathbf{q}_1 , the variance of the coordinates of $\mathbf{x}_1, \dots, \mathbf{x}_n$ is maximized by \mathbf{q}_2 , and so on.

THE EIGENVECTOR BASIS COORDINATES

The spectral decomposition of the covariance matrix $\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ is usually expressed as $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues and the columns of \mathbf{Q} are orthonormal unit eigenvectors.

THE EIGENVECTOR BASIS COORDINATES

The spectral decomposition of the covariance matrix $\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ is usually expressed as $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues and the columns of \mathbf{Q} are orthonormal unit eigenvectors.

The coordinates of \mathbf{x}_1 with respect to the eigenvectors of the covariance matrix are precisely the entries of

$$\mathbf{Q}^T\mathbf{x}_1 = \begin{bmatrix} - & \mathbf{q}_1 & - \\ & \vdots & \\ - & \mathbf{q}_d & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x}_1 \\ | \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1^T\mathbf{x}_1 \\ \vdots \\ \mathbf{q}_d^T\mathbf{x}_1 \end{bmatrix}.$$

THE EIGENVECTOR BASIS COORDINATES

The spectral decomposition of the covariance matrix $\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ is usually expressed as $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ where $\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues and the columns of \mathbf{Q} are orthonormal unit eigenvectors.

The coordinates of \mathbf{x}_1 with respect to the eigenvectors of the covariance matrix are precisely the entries of

$$\mathbf{Q}^T\mathbf{x}_1 = \begin{bmatrix} - & \mathbf{q}_1 & - \\ & \vdots & \\ - & \mathbf{q}_d & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x}_1 \\ | \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1^T\mathbf{x}_1 \\ \vdots \\ \mathbf{q}_d^T\mathbf{x}_1 \end{bmatrix}.$$

Indeed, one can easily compute the coordinates of all the observations at once:

$$\tilde{\mathbf{X}}\mathbf{Q} = \begin{bmatrix} - & \mathbf{x}_1 & - \\ & \vdots & \\ - & \mathbf{x}_n & - \end{bmatrix} \begin{bmatrix} | & & | \\ \mathbf{q}_1 & \cdots & \mathbf{q}_d \\ | & & | \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T\mathbf{q}_1 & \cdots & \mathbf{x}_1^T\mathbf{q}_d \\ \vdots & \ddots & \vdots \\ \mathbf{x}_n^T\mathbf{q}_1 & \cdots & \mathbf{x}_n^T\mathbf{q}_d \end{bmatrix}.$$

The i th row contains the coordinates of \mathbf{x}_i with respect to $\mathbf{q}_1, \dots, \mathbf{q}_n$.

THE PRINCIPAL COMPONENTS

An often more convenient matrix is $\tilde{\mathbf{X}}\mathbf{Q}$ which contains the coordinates of the *centered* observations with respect to $\mathbf{q}_1, \dots, \mathbf{q}_n$. In that case, the columns are orthogonal.

$$\begin{aligned}(\tilde{\mathbf{X}}\mathbf{Q})^T(\tilde{\mathbf{X}}\mathbf{Q}) &= \mathbf{Q}^T\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}\mathbf{Q} \\ &= \mathbf{Q}^T(n\mathbf{Q}\mathbf{A}\mathbf{Q}^T)\mathbf{Q} \\ &= n\mathbf{A}\end{aligned}$$

THE PRINCIPAL COMPONENTS

An often more convenient matrix is $\tilde{\mathbf{X}}\mathbf{Q}$ which contains the coordinates of the *centered* observations with respect to $\mathbf{q}_1, \dots, \mathbf{q}_n$. In that case, the columns are orthogonal.

$$\begin{aligned}(\tilde{\mathbf{X}}\mathbf{Q})^T(\tilde{\mathbf{X}}\mathbf{Q}) &= \mathbf{Q}^T\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}\mathbf{Q} \\ &= \mathbf{Q}^T(n\mathbf{Q}\mathbf{A}\mathbf{Q}^T)\mathbf{Q} \\ &= n\mathbf{A}\end{aligned}$$

These columns that were called *principal components* of the dataset, though the terminology is used inconsistently. You may see the columns of $\mathbf{X}\mathbf{Q}$ called the principal components. Or you may see $\mathbf{q}_1, \dots, \mathbf{q}_n$ called the principal components.

SINGULAR VALUE DECOMPOSITION

The eigenvectors and eigenvalues of the covariance matrix can also be obtained directly from a *singular value decomposition* of $\frac{1}{\sqrt{n}}\tilde{\mathbf{X}}$:

$$\frac{1}{\sqrt{n}}\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

where both \mathbf{U} and \mathbf{V} have orthonormal columns and \mathbf{D} is diagonal with non-negative entries.

SINGULAR VALUE DECOMPOSITION

The eigenvectors and eigenvalues of the covariance matrix can also be obtained directly from a *singular value decomposition* of $\frac{1}{\sqrt{n}}\tilde{\mathbf{X}}$:

$$\frac{1}{\sqrt{n}}\tilde{\mathbf{X}} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

where both \mathbf{U} and \mathbf{V} have orthonormal columns and \mathbf{D} is diagonal with non-negative entries.

How does this relate to the spectral decomposition of the covariance matrix?

$$\begin{aligned}\frac{1}{n}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}} &= \left(\frac{1}{\sqrt{n}}\tilde{\mathbf{X}}\right)^T\left(\frac{1}{\sqrt{n}}\tilde{\mathbf{X}}\right) \\ &= (\mathbf{U}\mathbf{D}\mathbf{V}^T)^T(\mathbf{U}\mathbf{D}\mathbf{V}^T) \\ &= \mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T \\ &= \mathbf{V}\mathbf{D}^2\mathbf{V}^T\end{aligned}$$

reveals a spectral decomposition of the covariance matrix. The columns of \mathbf{V} are unit eigenvectors while the squares of the diagonals in \mathbf{D} are the corresponding eigenvalues.

VARIANCES EQUAL TO EIGENVALUES

The variance of the coordinates with respect to \mathbf{q}_1 is

$$\begin{aligned}\mathbf{q}_1^T \left(\frac{1}{n} \widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}} \right) \mathbf{q}_1 &= \mathbf{q}_1^T (\lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \dots + \lambda_d \mathbf{q}_d \mathbf{q}_d^T) \mathbf{q}_1 \\ &= \lambda_1 \underbrace{(\mathbf{q}_1^T \mathbf{q}_1)^2}_1 + \dots + \lambda_d \underbrace{(\mathbf{q}_1^T \mathbf{q}_d)^2}_0 \\ &= \lambda_1.\end{aligned}$$

VARIANCES EQUAL TO EIGENVALUES

The variance of the coordinates with respect to \mathbf{q}_1 is

$$\begin{aligned}\mathbf{q}_1^T \left(\frac{1}{n} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right) \mathbf{q}_1 &= \mathbf{q}_1^T (\lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \dots + \lambda_d \mathbf{q}_d \mathbf{q}_d^T) \mathbf{q}_1 \\ &= \lambda_1 \underbrace{(\mathbf{q}_1^T \mathbf{q}_1)^2}_1 + \dots + \lambda_d \underbrace{(\mathbf{q}_1^T \mathbf{q}_d)^2}_0 \\ &= \lambda_1.\end{aligned}$$

Likewise the variance of the coordinates with respect to any eigenvector \mathbf{q}_j is the corresponding eigenvalue λ_j . If only the first $k \leq d$ principal components are retained, then the sum of the variances in the lower-dimensional approximation is $\lambda_1 + \dots + \lambda_k$.

VARIANCE DECOMPOSITION

The sum of the variances of the eigenvector basis coordinates equals the sum of the variances of the original coordinates.

VARIANCE DECOMPOSITION

The sum of the variances of the eigenvector basis coordinates equals the sum of the variances of the original coordinates.

A short proof is provided here for your future reference. It uses the fact that the squared norm of a vector is the sum of its squared coordinates with respect to any orthonormal basis (Parseval's identity). Without loss of generality, pretend once again that $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean $\mathbf{0}$.

$$\sum_{j=1}^d \frac{1}{n} \sum_{i=1}^n (\mathbf{q}_j^T \mathbf{x}_i)^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d (\mathbf{q}_j^T \mathbf{x}_i)^2 = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d x_{i,j}^2 = \sum_{j=1}^d \frac{1}{n} \sum_{i=1}^n x_{i,j}^2$$

VARIANCE DECOMPOSITION

The sum of the variances of the eigenvector basis coordinates equals the sum of the variances of the original coordinates.

A short proof is provided here for your future reference. It uses the fact that the squared norm of a vector is the sum of its squared coordinates with respect to any orthonormal basis (Parseval's identity). Without loss of generality, pretend once again that $\mathbf{x}_1, \dots, \mathbf{x}_n$ have mean $\mathbf{0}$.

$$\sum_{j=1}^d \frac{1}{n} \sum_{i=1}^n (\mathbf{q}_j^T \mathbf{x}_i)^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d (\mathbf{q}_j^T \mathbf{x}_i)^2 = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d x_{i,j}^2 = \sum_{j=1}^d \frac{1}{n} \sum_{i=1}^n x_{i,j}^2$$

Let v equal the sum of the variances of the original coordinates. If only the first $k \leq d$ principal components are retained, then the proportion of the variance that is preserved is $\frac{\lambda_1 + \dots + \lambda_k}{v}$.

SCALE-DEPENDENCE

Although variance doesn't depend on location, it does depend on *scale*. If a set of heights were measured in inches, they will have a much larger variances than they would if measured in feet. Likewise, the choice of units of measurement also affect which basis vectors are returned by the PCA algorithm.

SCALE-DEPENDENCE

Although variance doesn't depend on location, it does depend on *scale*. If a set of heights were measured in inches, they will have a much larger variances than they would if measured in feet. Likewise, the choice of units of measurement also affect which basis vectors are returned by the PCA algorithm.

To avoid having PCA depend on the units of measurement, the original variables can be rescaled to each have variance 1 before starting the algorithm.

SCALE-DEPENDENCE

Although variance doesn't depend on location, it does depend on *scale*. If a set of heights were measured in inches, they will have a much larger variances than they would if measured in feet. Likewise, the choice of units of measurement also affect which basis vectors are returned by the PCA algorithm.

To avoid having PCA depend on the units of measurement, the original variables can be rescaled to each have variance 1 before starting the algorithm.

However, if all of the original variables are measured in the same units, you may want to preserve their differing variances rather than rescaling.

IMPLEMENTATION IN R

RUNNING EXAMPLE CODE TOGETHER

See *PCA-examples.R* Exercises 1 - 3.

EXERCISES TO TRY FOR YOURSELF

Work through *PCA-examples.R* Exercise 4.