

Subspace Classifiers

Robert M. Haralick

Computer Science, Graduate Center
City University of New York

Outline

- 1 Projection Operators
- 2 Principal Components
- 3 Subspace Classifiers

Cartesian Products

Definition

The *Cartesian Product* of sets A_1, \dots, A_K is written as

$$A_1 \times A_2 \times \dots \times A_K$$

and is defined by

$$A_1 \times A_2 \times \dots \times A_K = \left\{ \left(\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_K \end{array} \right) \mid x_1 \in A_1, x_2 \in A_2, \dots, x_K \in A_K \right\}$$

The set A^K is called the *K-fold Cartesian Product* of A.

Euclidean Space

Definition

\mathbb{R} represents the set of all real numbers

Definition

An *N -dimensional Euclidean Space* is the set of all N -tuples of real numbers written as \mathbb{R}^N

Definition

The *Dimension* of \mathbb{R}^N is N

All the spaces we work with are Euclidean Spaces

Subspace Example

First we have to know what is a space or subspace

A three dimensional Euclidean Space has three kinds of subspaces:

- The zero dimensional point at the origin
- A one dimensional line
 - of infinite extent
 - of arbitrary orientation
 - containing the origin
- A two dimensional plane
 - of infinite extent
 - of arbitrary orientation
 - containing the origin

Scalars and Linear Spaces

Definition

A **Scalar** is any number from \mathbb{R}

Definition

A space \mathcal{L} is called a **Linear Space** if and only if for every scalar α and β

- $x \in \mathcal{L}$ and $y \in \mathcal{L}$ implies that $\alpha x + \beta y \in \mathcal{L}$

Any $x \in \mathcal{L}$ is called a point or vector of \mathcal{L}

Subspace

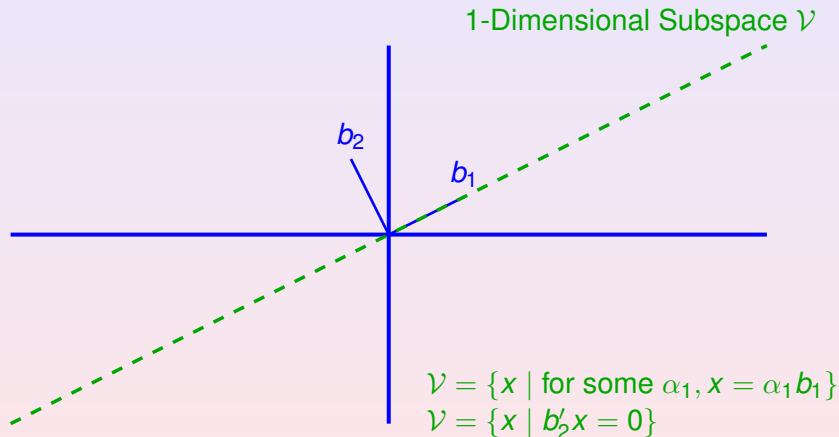
Definition

A subset $\mathcal{V} \subseteq \mathcal{L}$ is called a *Linear subspace* of \mathcal{L} if and only if for every scalars α and β

- $x \in \mathcal{V}$ and $y \in \mathcal{V}$ implies that $\alpha x + \beta y \in \mathcal{V}$

We are only interested in spaces and subspaces that are linear

Representing Subspaces



2-Dimensional Space

Linear Combination

Definition

A vector x from a subspace \mathcal{V} is said to be a *linear combination* of vectors x_1, \dots, x_K if and only if for some scalars $\alpha_1, \dots, \alpha_K$

$$x = \sum_{n=1}^N \alpha_n x_n \quad (1)$$

Linear Independence

Definition

A set of vectors $\{x_1, x_2, \dots, x_K\}$ from a subspace \mathcal{V} is said to be *Linearly Independent* if and only if for every set of scalars $\{\alpha_1, \dots, \alpha_K\}$, not all zero,

$$\sum_{k=1}^K \alpha_k x_k \neq 0$$

The vectors $\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$ are linearly independent

Linear Dependence

Definition

A set of vectors $\{x_1, \dots, x_K\}$ is said to be *Linearly Dependent* if and only if it is not linearly independent.

The vectors $\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}, \begin{pmatrix} 3 \\ 6 \\ 9 \end{pmatrix}$ are linearly dependent

Linear Dependence

Proposition

A set of vectors $\{x_1, x_2, \dots, x_K\}$ from a subspace V is Linearly Dependent if for some set of scalars $\{\alpha_1, \dots, \alpha_K\}$, not all zero,

$$\sum_{k=1}^K \alpha_k x_k = 0$$

Span

Definition

The *Span* of a set B ,

$$B = \{b_1, \dots, b_K \mid b_k \in \mathbb{R}^N, k = 1, \dots, K\}$$

is the set of all linear combinations of vectors from B

- We denote the span of B by $\text{Span}\{B\}$
- There is no constraint on K relative to N
- $\text{Span}\{B\}$ is a subspace of \mathbb{R}^N

Basis

Definition

A set B of vectors

$$B = \{b_1, \dots, b_K\}$$

is called a **Basis** for the subspace $\text{Span}\{B\}$ if and only if B is a linearly independent set

The vectors $a = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$, $b = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$ are linearly independent and constitute a basis for the subspace $\text{Span}\{a, b\}$

Dimension

Definition

The *Dimension* of a subspace \mathcal{V} is the smallest integer K such that the span of $\{b_1, \dots, b_K \mid b_k \in \mathcal{V}, k = 1, \dots, K\}$ satisfies

$$\text{Span}\{b_1, \dots, b_K\} = \mathcal{V}$$

Proposition

The dimension of a subspace \mathcal{V} is the largest number of vectors from \mathcal{V} such that the vectors are linearly independent.

Inner Product

Definition

The *Inner Product* between two vectors a and b from subspace \mathcal{V} is denoted by $a \cdot b$.

$$\text{Let } a = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_K \end{pmatrix} \text{ and } b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_K \end{pmatrix}$$

$$\text{Then, } a \cdot b = \sum_{k=1}^K a_k b_k$$

The dot product between vectors a and b can also be written in matrix notation

$$a \cdot b = a'b$$

Orthogonality

Definition

Two vectors $a, b \in \mathcal{V}$ are said to be orthogonal if and only if

$$a \cdot b = 0$$

We express that two vectors a and b are orthogonal by $a \perp b$.

Two spaces \mathcal{V} and \mathcal{W} are said to be orthogonal if and only if for every $v \in \mathcal{V}$ and every $w \in \mathcal{W}$

$$v \cdot w = 0$$

We express that two subspaces \mathcal{V} and \mathcal{W} are orthogonal by $\mathcal{V} \perp \mathcal{W}$

Direct Sum

Definition

Let \mathcal{V} and \mathcal{W} be subspaces of \mathcal{S} . The *Direct Sum* of subspaces \mathcal{V} and \mathcal{W} , denoted by $\mathcal{V} \oplus \mathcal{W}$, is defined by

$$\mathcal{V} \oplus \mathcal{W} = \{x \in \mathcal{S} \mid \text{for some } v \in \mathcal{V} \text{ and some } w \in \mathcal{W}, x = v + w\}$$

Orthogonal Complement

Definition

Let \mathcal{V} be a subspace of \mathcal{S} . The orthogonal complement of \mathcal{V} with respect to \mathcal{S} is denoted by \mathcal{V}^\perp and is defined by

$$\mathcal{V}^\perp = \{w \in \mathcal{S} \mid \text{for every } v \in \mathcal{V}, w \perp v\}$$

Orthogonal Complement Subspace

Definition

Let \mathcal{V} and \mathcal{W} be two subspaces of \mathcal{S} . \mathcal{W} is called the orthogonal complement of \mathcal{V} if and only if $\mathcal{W} = \mathcal{V}^\perp$

Proposition

Let \mathcal{V} be a subspace of \mathcal{S} . Then

- $\mathcal{V} \perp \mathcal{V}^\perp$
- $\mathcal{V} \oplus \mathcal{V}^\perp = \mathcal{S}$

Orthogonal Basis

Definition

A basis B for a subspace \mathcal{V} is said to be an orthogonal basis if and only for every $x, y \in B, x \neq y, x \perp y$

Orthogonal Projection

Definition

Let \mathcal{V} be a subspace of \mathcal{S} . Let $x \in \mathcal{S}$ and $x = v + w$ where $v \in \mathcal{V}$ and $w \in \mathcal{V}^\perp$. Then v is called the orthogonal projection of x onto \mathcal{V}

The orthogonal projection of x is unique.

Proposition

Let \mathcal{V} be a subspace of \mathcal{S} . Let $x \in \mathcal{S}$ and $x = v_1 + w_1 = v_2 + w_2$ where $v_1, v_2 \in \mathcal{V}$ and $w_1, w_2 \in \mathcal{V}^\perp$. Then $v_1 = v_2$

Proposition

Let \mathcal{V} be a subspace of S and $x \in S$. Let $x = v + w$ where $v \in \mathcal{V}$ $w \in \mathcal{V}^\perp$. Then $\|x\|^2 = \|v\|^2 + \|w\|^2$.

Proof.

Since $x = v + w$,

$$\begin{aligned}
 \|x\|^2 &= \|v + w\|^2 \\
 &= (v + w)'(v + w) \\
 &= v'v + v'w + w'v + w'w \\
 &= v'v + w'w \\
 &= \|v\|^2 + \|w\|^2
 \end{aligned}$$



Projection Operator

Definition

A square matrix P is said to be a **projection operator** if and only if

$$P^2 = P$$

A square matrix P is said to be an **orthogonal projection operator** if and only if

$$P^2 = P$$

$$P = P'$$

Projection Operator Examples

Definition

P is called a **projection operator** if and only if $P^2 = P$

$$\begin{pmatrix} .3 & .7 \\ .3 & .7 \end{pmatrix} \begin{pmatrix} .3 & .7 \\ .3 & .7 \end{pmatrix} = \begin{pmatrix} .3 & .7 \\ .3 & .7 \end{pmatrix}$$

$$\begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix}$$

$$\begin{pmatrix} .2 & .4 \\ .4 & .8 \end{pmatrix} \begin{pmatrix} .2 & .4 \\ .4 & .8 \end{pmatrix} = \begin{pmatrix} .2 & .4 \\ .4 & .8 \end{pmatrix}$$

Orthogonal Projection Operator Example

Consider the orthogonal projection operator onto the space spanned by

$$\frac{1}{5} \begin{pmatrix} 3 \\ 4 \end{pmatrix}$$

$$P = \frac{1}{5} \begin{pmatrix} 3 \\ 4 \end{pmatrix} \frac{1}{5} (3 \ 4) = \frac{1}{25} \begin{pmatrix} 9 & 12 \\ 12 & 16 \end{pmatrix}$$

$$\frac{1}{25} \begin{pmatrix} 9 & 12 \\ 12 & 16 \end{pmatrix} = \begin{pmatrix} \frac{3}{5} & \frac{-4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \frac{3}{5} & \frac{4}{5} \\ \frac{-4}{5} & \frac{3}{5} \end{pmatrix}$$

Orthogonal Projection Operator

Proposition

Let $\{b_1, b_2, \dots, b_K\}$ be an orthonormal basis for subspace \mathcal{V} , which is a subspace of an N -dimensional space S . Then the orthogonal projection operator P onto the K -dimensional subspace \mathcal{V} can be constructed by

$$P^{N \times N} = B^{N \times K} B'^{K \times N}$$

where

$$B = \begin{pmatrix} \vdots & \vdots & & \vdots \\ b_1 & b_2 & \dots & b_K \\ \vdots & \vdots & & \vdots \end{pmatrix}$$

Proposition

Let b_1, \dots, b_K be an orthonormal basis for the subspace \mathcal{V} . Let B be a matrix whose columns are the basis elements. Then, BB' is an orthogonal projection operator onto Col

Proof.

$$\begin{aligned}(BB')(BB') &= B(B'B)B^{-1}BB' \\ (BB')' &= BB'\end{aligned}$$

Finally, since any vector x that BB' operates on results in a linear combination of the columns of B , the space that BB' projects to is $\text{Col}(B)$ and the columns of B are the orthonormal basis vectors for the subspace \mathcal{V} . Hence BB' is the orthogonal projection operator onto the space \mathcal{V} . □

Orthogonal Projection Operators are Unique

Proposition

Let Q and P be orthogonal projection operators to the same subspace \mathcal{V} . Then $Q = P$

Proof.

Since both P and Q are orthogonal projection operators to the same subspace \mathcal{V} , the columns of P and the columns of Q lie in \mathcal{V} . Hence $PQ = Q$ and $QP = P$. Since Q is an orthogonal projection operator $Q = Q'$ and $PQ = Q$. Therefore,

$$Q = Q' = (PQ)' = Q'P' = QP = P$$



Another Form For Orthogonal Projection Operators

Proposition

Let b_1, \dots, b_K be an orthonormal basis for \mathcal{V} . Then $\sum_{k=1}^K b_k b_k'$ is an orthogonal projection operator onto the subspace \mathcal{V} .

Proof.

$$\begin{aligned} \sum_{j=1}^K b_j b_j' \sum_{k=1}^K b_k b_k' &= \sum_{j=1}^K \sum_{k=1}^K b_j (b_j' b_k) b_k' \\ &= \sum_{j=1}^K b_j b_j' \\ \left(\sum_{k=1}^K b_k b_k' \right)' &= \sum_{k=1}^K (b_k b_k')' = \sum_{k=1}^K b_k b_k' \end{aligned}$$

It is clear that whenever $\sum_{k=1}^K b_k b_k'$ operates on x , the result is a linear combination of the basis vectors for \mathcal{V} . □

Orthogonal Projection Minimizes Error

Theorem

Let \mathcal{V} be a subspace of S . Let $f : S \rightarrow \mathcal{V}$ and $x \in S$.

$$\min_f (x - f(x))' (x - f(x))$$

is achieved when f is the orthogonal projection operator from S to \mathcal{V}

Proof: Orthogonal Projection Minimizes Error

Proof.

Let $x \in \mathcal{S}$. Then there exists $v \in \mathcal{V}$ and $w \in \mathcal{V}^\perp$ such that $x = v + w$. Consider

$$\begin{aligned}\epsilon^2 &= \|x - f(x)\|^2 \\ &= (x - f(x))'(x - f(x)) \\ &= x'x - (v + w)'f(x) - f(x)'(v + w) + f(x)'f(x)\end{aligned}$$

But $f(x) \in \mathcal{V}$ and $w \in \mathcal{V}^\perp$. Hence $w'f(x) = 0$, therefore

$$\begin{aligned}\epsilon^2 &= x'x - v'f(x) + f(x)'v + f(x)'f(x) \\ &= (v + w)'(v + w) - v'f(x) - f(x)'v + f(x)'f(x) \\ &= v'v + w'w - v'f(x) - f(x)'v + f(x)'f(x) \\ &= (v - f(x))'(v - f(x)) + w'w\end{aligned}$$

ϵ^2 is minimized by making $f(x) = v$, the orthogonal projection of x into \mathcal{V} . □

Dimensional Reduction by Orthogonal Projection

Corollary

Let $x_1, \dots, x_K \in S$. Let \mathcal{V} be a subspace of S . Then

$$\min_{f: S \rightarrow \mathcal{V}} \sum_{k=1}^K \|x_k - f(x_k)\|^2$$

is achieved when f is the orthogonal projection operator from S to \mathcal{V}

Proof.

The best f can do for each x_k is for $f(x_k) = v_k$, the orthogonal projection of x_k onto \mathcal{V} . Therefore,

$$\min_{f: S \rightarrow \mathcal{V}} \sum_{k=1}^K (x_k - f(x_k))'(x_k - f(x_k))$$

is achieved when f is the orthogonal projection operator onto \mathcal{V} . □

Trace

Definition

The **Trace** of a $K \times K$ square matrix $A = (a_{ij})$ is defined by

$$\text{Trace}(A) = \sum_{k=1}^K a_{kk}$$

Proposition

Let A, B, A_1, \dots, A_K be square $N \times N$ matrices and $\alpha, \beta, \alpha_1, \dots, \alpha_K$ be scalars. Then

- $\text{Trace}(AB) = \text{Trace}(BA)$
- *Trace is a linear operator*
 - $\text{Trace}(\alpha A + \beta B) = \alpha \text{Trace}(A) + \beta \text{Trace}(B)$
 - $\text{Trace}(\sum_{k=1}^K \alpha_k A_k) = \sum_{k=1}^K \alpha_k \text{Trace}(A_k)$

Trace of Orthogonal Projection Operator

Proposition

Let P be an orthogonal projection operator to the M dimensional subspace \mathcal{V} . Then $\text{Trace}(P) = M$

Proof.

Let b_1, \dots, b_M be an orthonormal basis for \mathcal{V} . Then $P = \sum_{m=1}^M b_m b'_m$

$$\begin{aligned} \text{Trace}(P) &= \text{Trace}\left(\sum_{m=1}^M b_m b'_m\right) \\ &= \sum_{m=1}^M \text{Trace}(b_m b'_m) = \sum_{m=1}^M \text{Trace}(b'_m b_m) \\ &= \sum_{m=1}^M \text{Trace}(1) = \sum_{m=1}^M 1 = M \end{aligned}$$



Kernel and Range

Definition

- The **Kernel** of a matrix operator A is

$$\text{Kernel}(A) = \{x \mid Ax = 0\}$$

- The **Range** of a matrix operator A is

$$\text{Range}(A) = \{y \mid \text{for some } x, y = Ax\}$$

Definition

The **Column Space** of a matrix A is denoted by $\text{Col}(A)$ is defined by the Span of its columns.

Proposition

The Range(A) = Col(A)

Minkowski Sum

Definition

The **Minkowski Sum** or simply **Sum** of two subsets A and B of \mathcal{S} is defined by

$$A \oplus B = \{x \in \mathcal{S} \mid \text{for some } a \in A \text{ and for some } b \in B, x = a + b\}$$

Kernel and Range

Proposition

Let P be a projection operator onto subspace \mathcal{V} of S . Then

$$\text{Range}(P) \oplus \text{Ker}(P) = S$$

Proof.

Let $x \in S$. $Px + (I - P)x = Px + x - Px = x$. Certainly $Px \in \text{Range}(P)$. Consider $(I - P)x$.

$P[(I - P)x] = Px - PPx = Px - Px = 0$ Therefore, by definition of Kernel(P), $(I - P)x \in \text{Kernel}(P)$. □

Kernel and Range of Orthogonal Projection Operator

Proposition

Let P be an orthogonal projection operator. Then $\text{Range}(P) \perp \text{Kernel}(P)$

Proof.

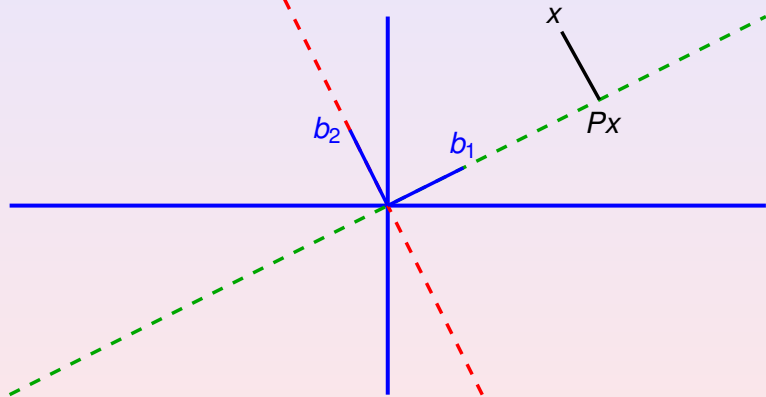
Let $x \in \text{Range}(P)$ and $y \in \text{Kernel}(P)$. Then for some u , $x = Pu$. Consider $x'y$.

$$x'y = (Pu)'y = u'P'y = u'Py$$

But $y \in \text{Kernel}(P)$ so that $Py=0$. Therefore $x'y = 0$. □

Range And Kernel of Orthogonal Projection Operator

$$P(\alpha b_2) = b_1 b_1' (\alpha b_2) = 0 \quad P = b_1 b_1' \quad P(\alpha b_1) = b_1 b_1' (\alpha b_1) = \alpha b_1$$



$$\text{Range}(P) = \{x \in \mathcal{S} \mid \text{for some } y, x = P(y)\}$$

$$\text{Kernel}(P) = \{x \in \mathcal{S} \mid Px = 0\}$$

2-Dimensional Space \mathcal{S}

Projection Operator to \mathcal{V}^\perp

Proposition

Let P be a projection operator onto the subspace \mathcal{V} . (Not necessarily an orthogonal projection operator) Then $I - P$ is the projection operator onto the subspace \mathcal{V}^\perp .

Proof.

$$\begin{aligned}(I - P)(I - P) &= I - P - P + P^2 \\ &= I - 2P + P = I - P\end{aligned}$$

$I - P$ is also a projection operator. But what space does it project to?

$\mathcal{V}^\perp = \text{Kernel}(P)$. Let $x \in \mathcal{V}^\perp$. Then $Px = 0$. Consider $(I - P)x = x - Px = x$



Orthogonal Projection Operator to \mathcal{V}^\perp

Proposition

Let P be the orthogonal projection operator onto the subspace \mathcal{V} . Then $I - P$ is the orthogonal projection operator onto the subspace \mathcal{V}^\perp .

Proof.

We already know that $(I - P)(I - P) = I - P$. We just have to show that $I - P$ is symmetric and that $I - P$ projects to the \mathcal{V}^\perp .

$$(I - P)' = I' - P' = I - P$$

$\mathcal{V}^\perp = \text{Kernel}(P)$. Let $x \in \mathcal{V}^\perp$. Then $Px = 0$. Consider $(I - P)x = x - Px = x$ Every $x \in \mathcal{V}^\perp$ gets projected to itself. \square

Covariance Matrix and Expected Value of Squared Length

Definition

The **Covariance Matrix** of a distribution is defined by

$$\Sigma = E[(x - \mu)(x - \mu)']$$

Proposition

$$\text{Trace}(\Sigma) = E[\|x - \mu\|^2]$$

Proof.

$$\begin{aligned} \text{Trace}(\Sigma) &= \text{Trace}(E[(x - \mu)(x - \mu)']) \\ &= E[\text{Trace}(x - \mu)(x - \mu)'] \\ &= E[\text{Trace}((x - \mu)'(x - \mu))] \\ &= E[(x - \mu)'(x - \mu)] \\ &= E[\|x - \mu\|^2] \end{aligned}$$

Covariance Matrix and Sum of Squared Vector Lengths

Given a sample x_1, \dots, x_M of N -dimensional vectors, the unbiased estimated of the covariance matrix Σ is given by

$$\Sigma = \frac{1}{M-1} \sum_{m=1}^M (x_m - \mu)(x_m - \mu)'$$

where the estimated mean μ is given by

$$\mu = \frac{1}{M} \sum_{m=1}^M x_m$$

Then

$$\text{Trace}(\Sigma) = \frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2$$

Covariance and Sum of Squared Vector Lengths

Proposition

Let Σ be the unbiased estimated covariance matrix. Then

$$\text{Trace}(\Sigma) = \frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2$$

Proof.

$$\begin{aligned} \text{Trace}(\Sigma) &= \text{Trace} \left(\frac{1}{M-1} \sum_{m=1}^M (x_m - \mu)(x_m - \mu)' \right) = \frac{1}{M-1} \text{Trace} \left(\sum_{m=1}^M (x_m - \mu)(x_m - \mu)' \right) \\ &= \frac{1}{M-1} \sum_{m=1}^M \text{Trace} \left((x_m - \mu)(x_m - \mu)' \right) = \frac{1}{M-1} \sum_{m=1}^M \text{Trace} \left((x_m - \mu)'(x_m - \mu) \right) \\ &= \frac{1}{M-1} \sum_{m=1}^M (x_m - \mu)'(x_m - \mu) = \frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2 \end{aligned}$$

□

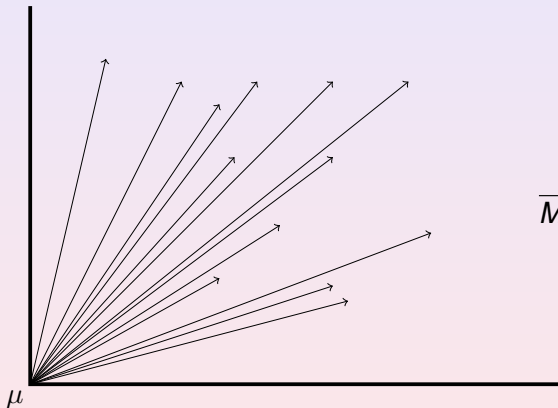
Covariance and Sum of Squared Vector Lengths

The covariance matrix Σ gives information about the spread of the vectors composing it. $\text{Trace}(\Sigma)$ is a measure of the total variance. The sum

$$\frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2$$

is a normalized sum of the squared length of the x_m vectors to the mean μ .

Sum of Squared Vector Lengths



$$\frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2$$

Eigen Decomposition of Covariance Matrix

Proposition

Let Σ be the covariance matrix. Let the Eigenvalue Eigenvector decomposition of Σ be $\Sigma = T\Lambda T'$ where $\Lambda = \text{Diagonal}(\lambda_1, \dots, \lambda_N)$
 Then,

$$\text{Trace}(\Sigma) = \sum_{n=1}^N \lambda_n$$

Proof.

$$\begin{aligned} \text{Trace}(\Sigma) &= \text{Trace}(T\Lambda T') = \text{Trace}(\Lambda T T') \\ &= \text{Trace}(\Lambda) \\ &= \sum_{n=1}^N \lambda_n \end{aligned}$$

Eigen Decomposition of Covariance Matrix

Proposition

Let Σ be the covariance matrix of random vector x . Let the Eigenvalue Eigenvector decomposition of Σ be $\Sigma = T\Lambda T'$ where $\Lambda = \text{Diagonal}(\lambda_1, \dots, \lambda_N)$. Let the n^{th} column of T be t_n . Let $\sigma_n^2 = V[t_n'x]$. Then,

$$\sigma_n^2 = \lambda_n$$

Proof.

$$\begin{aligned} \sigma_n^2 &= V[t_n'x] = E[(t_n'x - E[t_n'x])^2] \\ &= E[(t_n'(x - \mu))((x - \mu)'t_n)] \\ &= t_n'E[(x - \mu)(x - \mu)']t_n = t_n'\Sigma t_n \\ &= t_n'T\Lambda T't_n = (0 \dots 010 \dots 0)\Lambda(0 \dots 010 \dots 0)' \\ &= \lambda_n \end{aligned}$$

But if Σ and μ are estimated from the data,

$$\text{Trace}(\Sigma) = \frac{1}{M-1} \sum_{n=1}^N \|x_m - \mu\|^2$$

Therefore, we can conclude that

$$\sum_{n=1}^N \sigma_n^2 = \frac{1}{M-1} \sum_{m=1}^M \|x_m - \mu\|^2$$

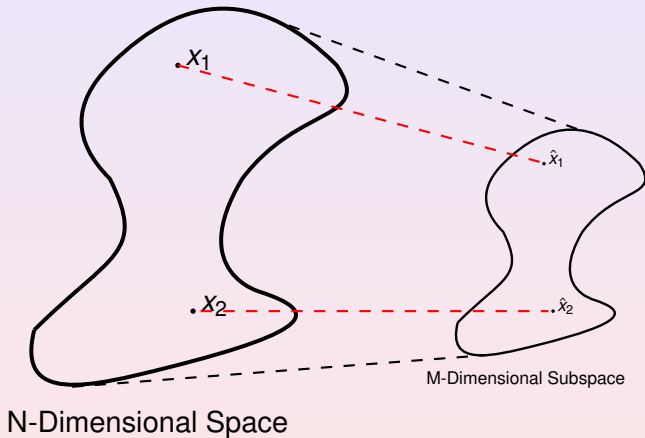
The sum of the Eigenvalues of the Covariance Matrix is the total variance and is equal to the $\frac{1}{M-1}$ of the sum of the squared length of the $(x_m - \mu)$ vectors.

Projecting to Subspaces

Principal Components projects the original data from the larger dimensional space in which it resides to a smaller dimensional space.

- If the decision is to project to a subspace of dimension K , which subspace should be chosen?
- With Principal Components, the K -dimensional space is found that minimizes the sum of the squared distances between the original data vectors and their projection in the K -dimensional subspace.

Space Squeezing: Dimensionality Reduction



Principal Components and Orthogonal Projection Operators

Consider the case for an orthogonal projection operator. It projects a data point or vector to that place in the subspace that is closest to the original point.

Suppose the original data points are N -dimensional. The projection operator projects each point to the closest point to it in the K -dimensional subspace determined by the range of the orthogonal projection operator.

For some subspaces of dimension K the overall distances between the original data points and their projections will be the smallest. This is the one that Principal Components determines.

The Simplest Orthonormal Projection Operators

The simplest orthogonal projection operator is a diagonal matrix with some of the entries on the diagonal being 1's and the other entries on the diagonal being 0's.

For example, examine the orthogonal projection operator that projects the data to its first two components, all other components of the projected vector being 0.

$$P = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix}$$

Clearly $P = P^2$ and $P = P'$ making it an orthogonal projection operator.

The Orthonormal Matrix

Definition

A square matrix Q is said to be *Orthonormal* if its columns each have norm 1 and each column is orthogonal to every other column.

Proposition

The transpose of an orthonormal matrix is its inverse.

Proof.

Let T be an $N \times N$ orthonormal matrix with columns t_1, \dots, t_N . Then

$$T' T = \begin{pmatrix} \dots & t'_1 & \dots \\ \dots & t'_2 & \dots \\ \vdots & \vdots & \vdots \\ \dots & t'_N & \dots \end{pmatrix} \begin{pmatrix} \vdots & \vdots & \dots & \vdots \\ t_1 & t_2 & \dots & t_N \\ \vdots & \vdots & \dots & \vdots \end{pmatrix} = I$$

$T' T = I$ so that $T' T T^{-1} = T^{-1}$. Hence,

$$T' = T^{-1}$$

The General Orthogonal Projection Operator

Proposition

If Q is an orthonormal matrix and P is an orthogonal projection operator, then QPQ' is an orthogonal projection operator.

Proof.

We have to show that QPQ' is idempotent and symmetric. Consider

$$\begin{aligned}(QPQ')(QPQ') &= QP(Q'Q)PQ' \\ &= QPPQ' \\ &= QPQ' \\ (QPQ')' &= QP'Q' \\ &= QPQ'\end{aligned}$$

The General Orthogonal Projection Operator

Proposition

Let P be an orthonormal projection operator. Let Q be an orthonormal matrix. Then QPQ' projects to a subspace of the same dimension as P

Proof.

Since the dimension of the space an orthonormal projection operator projects to is the trace of the operator, we just have to show that $\text{Trace}(P) = \text{Trace}(QPQ')$

$$\begin{aligned}\text{Trace}(QPQ') &= \text{Trace}(PQQ') \\ &= \text{Trace}(P(QQ')) \\ &= \text{Trace}(P)\end{aligned}$$



The Orthogonal Projection Operator In Diagonalized Form

Proposition

The form QPQ' can orthogonally project to any given subspace \mathcal{V} with P being a diagonal matrix have ones and zeros on the diagonal.

Proof.

Without loss of generality, we take $P^{N \times N}$ to be a diagonal matrix with the first $M < N$ entries being ones and the remaining diagonal entries zero. The proof is by construction. Let q_1, \dots, q_M be an orthonormal basis for \mathcal{V} . Extend this orthonormal basis to q_{M+1}, \dots, q_N . Define the matrix Q to have columns of q_1, \dots, q_N . Define the orthogonal projection operator P to be a diagonal matrix whose first M diagonal entries are one and all the remaining diagonal entries are zero. \square

Proof Continued

Consider QPQ' .

$$\begin{aligned}QPQ' &= \begin{pmatrix} \vdots & \vdots & \vdots \\ q_1 & \cdots & q_M \\ \vdots & \vdots & \vdots \end{pmatrix} \begin{pmatrix} \cdots & q'_1 & \cdots \\ \cdots & q'_2 & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & q'_M & \cdots \end{pmatrix} \\ &= \sum_{m=1}^M q_m q'_m\end{aligned}$$

And this is the orthogonal projection operator onto the subspace \mathcal{V}

The Principal Component Technique

Let x_1, \dots, x_K be the observed $N \times 1$ data vectors. First center the data around the mean by subtracting the sample mean vector μ from each of the original data points.

$$\mu = \frac{1}{K} \sum_{k=1}^K x_k$$

Define the sample unbiased covariance matrix Σ by

$$\Sigma = \frac{1}{K-1} \sum_{k=1}^K (x_k - \mu)(x_k - \mu)'$$

Eigenvector Eigenvalue Decomposition

Σ is an $N \times N$ real symmetric positive semidefinite matrix. Consider the eigenvalue eigenvector decomposition of Σ

$$\Sigma = U\Lambda U'$$

where Λ is a diagonal matrix of eigenvalues and U is an orthonormal matrix. Since Σ is a real symmetric positive semidefinite matrix the eigenvalues are non-negative.

Total Variance

The total variance is given by the trace of Σ . It has the meaning that it is $\frac{1}{K-1}$ times the squared distance between the observed data to the centroid given by the mean.

Note that the trace of Σ is equal to the trace of Λ

$$\begin{aligned} \text{Trace}(\Sigma) &= \text{Trace}(U\Lambda U') \\ &= \text{Trace}(\Lambda U U') \\ &= \text{Trace}(\Lambda) \end{aligned}$$

Total Variance

Without loss of generality, we suppose that the diagonal entries are ordered from largest to smallest. Since the eigenvalues are ordered in descending order, the first column of U is that subspace that would have the smallest distance between the observations and the mean vector in the subspace defined by the span of the first eigenvector. Alternatively it is also the subspace whose squared projected lengths is maximal.

The span of the second column of U would be that subspace, orthogonal to the first having the next most smallest squared distance between the observations and the mean vector. And so on.

Because the total variance is fixed, the sum of the squared distances between the data points and their projection are minimized.

Principal Components

Theorem

Let $x_1, \dots, x_K \in S$ an N -dimensional vector space and Q be an orthogonal projection operator of rank M . Then $\sum_{k=1}^K x_k Q x_k$ is maximized when Q projects to the M -Dimensional subspace spanned by the M eigenvectors of $\sum_{k=1}^K x_k x_k'$ having largest eigenvalues.

Proof.

Let $\sum_{k=1}^K x_k x_k' = T D T'$ and $Q^* = T' Q T$. Without loss of generality we assume that the diagonal entries are ordered $d_{ii} \geq d_{jj}$, $i < j$. Then $\max_{Q^*} \text{Trace}(Q^* D) = \sum_{m=1}^M d_{mm}$, where the maximum is taken over all Q^* satisfying $Q^* = Q^* Q^*$ and $Q^* = Q^{*'}'$. Thus, the first M diagonal entries of Q^* are one and the remaining diagonal entries 0. Since $\sum_{i=1}^N \sum_{j=1}^N q_{ij}^2 = M$, and there are M ones on the diagonal, the remaining elements of Q^* are 0. This implies $Q = T Q^* T'$ is the orthogonal projection operator onto the space spanned by the first M eigenvectors of $\sum_{k=1}^K x_k x_k'$ for these are the eigenvectors having largest eigenvalues. □

Choice of Best Subspace

In Principal Components the researcher calculates the successive sums of the eigenvalues and compares them to the total sum of the eigenvalues, which is the $Trace(\Lambda)$, and then sets the threshold. Calculate the running sum

$$v_n = \sum_{i=1}^n \lambda_i$$

Choice of Best Subspace

Choose the smallest n such that $\frac{v_n}{V_N}$ just exceeds the selected threshold θ . For example the threshold could be set to .85. n is chosen to be the smallest value satisfying

$$\frac{v_n}{V_N} > \theta$$

The subspace projected to is spanned by the first n columns of U . Suppose these n -columns are u_1, \dots, u_n . Then the orthogonal projection operator P is defined by

$$P = \sum_{i=1}^n u_i u_i'$$

Relative Coordinates

- Suppose every data point is an N -dimensional measurement from space \mathcal{S}
- Let $P^{N \times N}$ be a projection operator to M -dimensional subspace $\mathcal{V} \subset \mathcal{S}$
- Suppose b_1, \dots, b_M is any orthonormal basis for \mathcal{V}
- The projection operator P is given by

$$P^{N \times N} = \sum_{m=1}^M b_m b_m'$$

- $y^{N \times 1} = P^{N \times N} x^{N \times 1}$ is the projection of x into \mathcal{V}

Relative Coordinates

- Although y lies in a M -dimensional subspace \mathcal{V} , y is an N -dimensional vector
- Since $y \in \mathcal{V}$, we can write $y = \sum_{n=1}^M \alpha_n b_n$ since b_1, \dots, b_M is a basis for \mathcal{V}
- The tuple $(\alpha_1, \dots, \alpha_M)$ is called the relative coordinates of the projection of x
- Let $B^{N \times M}$ be a matrix whose M columns are the basis vectors b_1, \dots, b_M
- The coefficients can be obtained by $(\alpha_1, \dots, \alpha_M)' = B'x$
- Then the following calculation can produce the orthogonal projection y

$$y^{N \times 1} = B^{N \times M} (\alpha_1, \dots, \alpha_M)'^{M \times 1} = BB'x = Px$$

Principal Components

- Disregarding the class labels, Principle Components selects that K-dimensional subspace having the best fit to the observed measurement vectors
- For each measurement vector, Principal Components computes its relative coordinates in the subspace
- Classification is done using the relative coordinates

Feature Selection

- Feature Selection is the oldest form of subspace classifiers
- There are many papers describing ways of doing feature selection
- From one point of view, the problem of Feature Selection is to select a fixed number of features that will maximize the classification accuracy.
- The problem of selecting the best K of N features for the classification task is NP-Hard
- The three oldest techniques are
 - Forward Search
 - Backward Search
 - Combinatorial Search

Forward Feature Selection Greedy Approach

- Check every feature one at time
- Construct a classifier using the candidate feature
- Select the candidate feature with the highest classification accuracy
- Check every one of the remaining features
- Construct a classifier using the previously selected features and the current candidate selection
- Stop when there are no more improvements

Backward Feature Selection Greedy Approach

- Determine the classification accuracy using all the features
- Then check the classification accuracy by leaving out a feature
- Select that feature to leave out which decreases the classification the least
- Iterate checking what happens when you leave out a feature from the current set of features
- Continue the process until the classification accuracy decreases too much.

Combinatorial Feature Selection

- If there are N features there are $2^N - 1$ nonempty sets of features
- Go through all combinations of subsets of features
- For each combination determine the classification accuracy
- Choose that combination whose classification accuracy is highest

Feature Selection Is An Application of A Projection Operator

- Whatever the feature selection methodology
- It results in a selection of a subset of features
- In essence, it does a relative coordinate orthogonal projection to a subset of features
- And leaves out the remaining features

Non-relevant Features

- Suppose that one or more of the features have nothing to do with class c
- They provide no information relative to class c
- The range of these useless features can be large
- Using Principle Components cannot help
 - Useless features with a large range will dominate
 - Principle Components would then include it

Alternatives

- N Features indexed by $(1, \dots, N)$
- Go through all features
- Go through all classes
- With respect to each class, determine the classification accuracy using each feature alone
- This results in a table whose rows are the classes and whose columns are the features.
- If for some feature, the classification accuracy for class c is too small, mark the feature as nonessential for class c
- Compute the covariance matrix Σ_c for class c using only essential features

Details

- Let there be N features and K classes
- Take the features one by one and build a classifier using the Training Sequence
- Use the Test Sequence to determine how well the classifier does with each feature for each class
- Let $d(f, x)$ be the class that the classifier assigns to x when only using feature f

Let $\langle x_1, \dots, x_Z \rangle$ be the Test Sequence of N -tuples, with corresponding true class tags $\langle c_1, \dots, c_Z \rangle$ then the resulting N row by K column table A is defined by

$$A(f, c) = \frac{|\{z \in [1, Z] \mid d(f, x_z) = c \text{ and } c_z = c \text{ when feature } f \text{ is used}\}|}{|\{z \in [1, Z] \mid c_z = c\}|}$$

Given class c , $A(f, c)$ is the probability of correct identification when using feature f

Details

- Let θ be the threshold that determines when classification accuracy is large enough
- If for any class c and feature f , $A(f, c) > \theta$, then feature f can be used in playing a role in the classification for class c .

Project To Essential Features

- Suppose there are 6 features
- Features 1,4, and 5 are essential; $Q = \langle 1, 4, 5 \rangle$
- Relative Coordinate Project to components 1,4 and 5
- The relative coordinate projection of tuple $x = (x_1, x_2, x_3, x_4, x_5, x_6)$ is (x_1, x_4, x_5)
- We can write the projection in a general way
 - Let Q be the index list of essential features
 - In the example above $Q = \langle 1, 4, 5 \rangle$
 - The indices are ordered in ascending order
 - The relative coordinate projection of tuple (x_1, \dots, x_N) onto the essential features specified in Q is given by
 - $\pi_Q(x_1, \dots, x_N) = (x_i : i \in Q)$

Finding Subspaces

- For each class c
- Use the relative coordinate projection of the training set to determine the class covariance matrix Σ_c
- Use Σ_c to do a Principle Components
- E_c is selected so that the eigenvalue fraction

$$\frac{\sum_{n=1}^{E_c} \sigma_n}{\sum_{n=1}^N \sigma_n}$$

is just greater than the user specified fraction f

- Define the class subspace to be the span of the first E_c eigenvectors of Σ_c

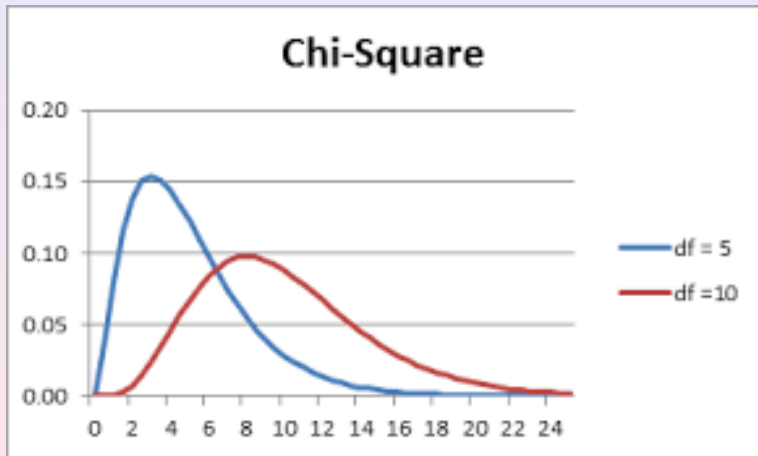
Known Mean and Covariance Matrix

- If the distribution were multivariate normal
- With known mean μ and known covariance matrix Σ
- The Mahalanobis distance of x to μ is given by
 - $d_c^2 = (x - \mu)' \Sigma^{-1} (x - \mu)$
- d_c^2 has a χ^2 distribution with E_c degrees of freedom
- If the mean is estimated from data with a known covariance matrix
- d_c^2 has a χ^2 distribution with $E_c - 1$ degrees of freedom

The χ^2 Distribution

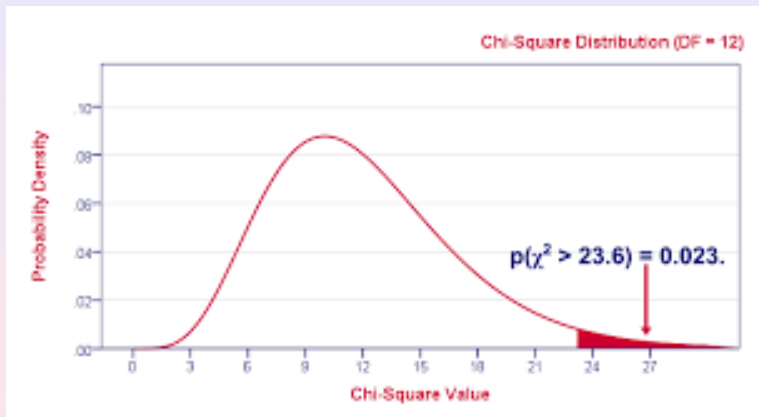
- The mean of a Chi Square distribution is its degrees of freedom
- Chi Square distributions are positively skewed (skewed to the right)
- Degree of skew decreases with increasing degrees of freedom
- As the degrees of freedom increases, the Chi Square distribution approaches a normal distribution
- Density function: $\frac{1}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} x^{\frac{k}{2}-1} e^{-\frac{x}{2}}$

The χ^2 Distribution



For critical value 16.75, the $Prob(\chi_5^2 > 16.75) = .005$

The χ^2 Distribution



For critical value 23.6, the $Prob(\chi_{25}^2 > 23.6) = 0.023$

Making The Class Assignment

- If the distribution were multivariate normal
- With known mean and known covariance matrix
- d_c^2 has a χ^2 distribution with E_c degrees of freedom
- If the mean is estimated from data with a known covariance matrix d_c^2 has a χ^2 distribution with $E_c - 1$ degrees of freedom
- $Prob(\chi_{E_c}^2 > s_{critical,c}(p_0)) = p_0$, $p_0 = .005$ or $.01$
- Consider assigning x to only those classes in S :

$$S = \{c \in C \mid d_c^2 \leq s_{critical,c}(p_0)\}$$
- Assign x to that class $c \in S$ with the smallest squared Mahalanobis distance

Making The Class Assignment

- For each class c
- Determine T_c , the matrix with orthonormal columns defined by the first E_c eigenvectors of Σ_c
- Choose p_0
- Define $s_{critical,c}(p_0)$ by $Prob(\chi_{E_c-1}^2 > s_{critical,c}(p_0)) = p_0$
- For each x in training set
- Using the relative coordinates of the subspace
 - The columns of $S_c^{N \times E_c}$ are defined by the first E_c eigenvectors of Σ_c
 - Do an orthogonal projection to the subspace associated with the class c : $y = S_c'(x - \mu)$
 - y has covariance matrix $S_c'\Sigma_c S_c$ which is a $E_c \times E_c$ matrix

Covariance Matrix of $y^{E_c \times 1}$

- $y^{E_c \times 1}$
- We need its covariance matrix so that we can use its inverse in the Mahalanobis distance calculation

$$\begin{aligned}
 \Sigma_y &= S'_c \Sigma S_c \\
 &= S'_c (T_c \Lambda T'_c) S_c \\
 &= (S'_c T_c) \Lambda (T'_c S_c) \\
 &= \begin{pmatrix} I^{E_c \times E_c} & 0^{E_c \times N - E_c} \end{pmatrix} \Lambda^{N \times N} \begin{pmatrix} I^{E_c \times E_c} \\ 0^{N - E_c \times E_c} \end{pmatrix} \\
 &= \text{Diagonal}(\lambda_1, \lambda_2, \dots, \lambda_{E_c})
 \end{aligned}$$

The inverse covariance matrix Σ_y^{-1}

$$\Sigma_y^{-1} = \text{Diagonal}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_{E_c}^{-1})$$

Mahalanobis Distance for y

$$d_c^2(y) = y' \text{Diagonal}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_{E_c}^{-1}) y$$

- If the distribution were multivariate normal
- $d_c^2(y)$ has $\chi_{E_c}^2$
- If $d_c^2(y) > s_{critical,c}(p_0)$, x cannot be assigned to class c
- Assign x to allowable class c for which $d_c^2(y)$ is minimal

What if x is not Multivariate Normal

- $y = S'_c x$
- Training sequence for class c is $\langle y_1, \dots, y_{Z_c} \rangle$
- Corresponding list of Squared Mahalanobis distances
 - $\langle d_c^2(y_1), \dots, d_c^2(y_{Z_c}) \rangle$
- Order them in ascending order
 - $\langle d_{(1)}^2, \dots, d_{(Z_c)}^2 \rangle$
- Choose $p_0 = \frac{Z}{Z_c}$
 - $s_{critical,c}(p_0) = d_{Z_c - Z}^2$
- If $d_c^2(y) > s_{critical,c}(p_0)$, class c is not allowable for x
- Assign x to allowable class c for which $d_c^2(y)$ is minimal

Review 1

- For class c
- Leave out components of x whose associated classification accuracy is too low
- The covariance matrix of the reduced x is Σ_c
- Eigenvector Eigenvalue decomposition $\Sigma_c = T_c \Lambda_c T_c'$
- Choose a fraction f of the variance to be preserved
- E_c is the smallest number satisfying $\frac{\sum_{n=1}^{E_c} \lambda_{cn}}{\sum_{n=1}^N \lambda_{cn}} \geq f$
- Define S_c to be the first E_c columns of T_c
- $y = S_c' x$
- y has covariance matrix $S_c' \Sigma_c S_c = \text{Diagonal}(\lambda_{c1}, \dots, \lambda_{cE_c})$
- The inverse covariance matrix is $\text{Diagonal}(\lambda_{c1}^{-1}, \dots, \lambda_{cE_c}^{-1})$

Review 2

- Squared Mahalanobis distance
 - $d_c^2(y) = y' \text{Diagonal}(\lambda_{c1}^{-1}, \dots, \lambda_{cE_c}^{-1})y$
- Training sequence for class c
 - $\langle x_1, \dots, x_{Z_c} \rangle$
- $y = S'_c x$
 - $\langle y_1, \dots, y_{Z_c} \rangle$
 - $\langle d_c^2(y_1), \dots, d_c^2(y_{Z_c}) \rangle$
 - Ascending order $\langle d_{(1)}^2, \dots, d_{(Z_c)}^2 \rangle$
 - $p_0 = \frac{z}{Z_c}$
 - $s_{critical,c}(p_0) = d_{(Z_c-z)}^2$
- New $x, y = S'_c x$
- If $d_c^2(y) > s_{critical,c}$ class c is not allowable for x
- Assign x to allowable class c for which $d_c^2(y)$ is minimal

Making The Class Assignment

- If the distribution were multivariate normal
- With known mean and known covariance matrix
- d_c^2 has a χ^2 distribution with E_c degrees of freedom
- If the mean is estimated from data with a known covariance matrix d_c^2 has a χ^2 distribution with $E_c - 1$ degrees of freedom
- Assign x to the class c with the smallest squared Mahalanobis distance d_c^2 , providing that $d_c^2 < s_{tail,c}$

Problem With Using The Mahalanobis Distance P-value

- It does not include the possibility maximizing economic gain
- Maximizing economic gain is easy with the Discrete Bayes Rule
- The Mahalanobis Distance P-value
 - Produces a real value between 0 and 1
 - The real value has to be converted to an integer to address the class conditional probability table
 - Solution is to quantize the Mahalanobis p-value for each class
- Quantizing
 - Equal Interval Quantizing
 - Equal Probability Quantizing

Equal Interval Quantizing

- Suppose we want K quantizing intervals
- The interval $[0, 1]$ is divided in equal subintervals of size $\frac{1}{K}$
- The quantizing boundaries are $\langle 0, \frac{1}{K}, \frac{2}{K}, \dots, \frac{K-1}{K}, 1 \rangle$
- Let p be a p-value
- If $\frac{k}{K} \leq p < \frac{k+1}{K}$ the quantizing index is k
- If $\frac{K-1}{K} \leq p \leq 1$ the quantizing index is K

Equal Probability Quantizing

- Suppose we want K quantizing intervals
- The indexes of the quantizing intervals range are in the set $\{0, 1, \dots, K - 1\}$
- The Training Sequence has Z tuples
- Z is a multiple of K : for some natural integer m , $Z = mK$
- Order the p-values in ascending order $p_{(1)}, p_{(2)}, \dots, p_{(Z)}$
- The quantizing boundaries are $\langle b_0 = 0, b_1, b_2, \dots, b_{K-1}, b_K = 1 \rangle$
- Where $b_k = p_{(kZ/K)}$, $k \in \{1, \dots, K - 1\}$
- If for some $k \in \{0, \dots, K - 1\}$, $b_k \leq p < b_{k+1}$, the quantizing index is k
- If $p \geq p_{(K-1)}$ the quantizing index is $K - 1$

Non-uniform Equal Probability Quantization

