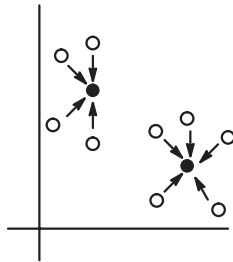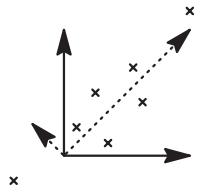# DATA COMPRESSION

There are two main ways of compressing numeric multidimensional data:



* Dimensionality reduction:

  + find directions that show maximal variation using *eigenvalue* techniques and neglect other directions.
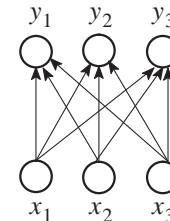
* Clustering or vector quantization: use *prototypes* to encode a group of points.

---

# COORDINATE TRANSFORMATIONS

* Linear transformation: $y = Wx$.
* Note that a single layer neural network without activation function also performs a linear transformation:

$y_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3$, etc.



* Suppose that a coordinate transformation is given by a matrix $A$:

$$x^* = Ax$$
$$y^* = Ay.$$

* Because $x = A^{-1}x^*$ it follows: $y^* = AWA^{-1}x^*$ or $y^* = W^*x^*$, with $W^* = AWA^{-1}$. $W^*$ and $W$ are called *similar* matrices.

---

# EIGENVALUES AND EIGENVECTORS

* For some vectors: $Wv = \lambda v$.

* Such a vector $v$ is called an *eigenvector* of $W$ and $\lambda$ is the corresponding *eigenvalue.*

* Consider the matrix $Y$ the columns of which are eigenvectors of $W$. Then: $WY = Y\Lambda$, with $\Lambda$ a *diagonal matrix*.

* Therefore: $Y^{-1}WY = Y^{-1}Y\Lambda = \Lambda$. This means that for each transformation $W$, there is an equivalent transformation by means of a diagonal matrix that uses the eigenvectors as a basis.

---

# RANDOM VECTORS

* A *random vector* is a vector whose components are random variables.

* A random vector $X$ has a probability density function $p(X)$.

* The *mean vector* is defined as: $M = E[X] = \int Xp(X)dX$ (integrate separately for each vector element).

* And the *covariance matrix* is defined as: $\Sigma = E\left[(X - M)(X - M)^T\right]$.

## SAMPLED RANDOM VECTORS

* In practice, the probability density function is unknown and one only has samples $X^k$ ($k = 1, ..., N$).
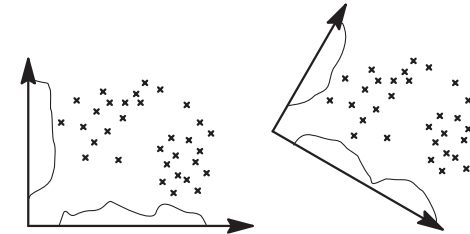
* The *sample mean vector* is defined as:

$$M = \frac{1}{N} \sum_{k=1}^{N} X^k.$$

* The *sample covariance matrix* is defined as:

$$\Sigma = \frac{1}{N} \sum_{k=1}^{N} (X^k - M)(X^k - M)^T$$

## PRINCIPAL COMPONENTS ANALYSIS (1)

* Abbreviated by PCA; also called *Karhunen-Loève transformation*.

* Question: given a random vector (a set of sampled vectors), find an orthogonal base that maximizes the variance along the subsequent dimensions of the base.

* Goal: achieve dimensionality reduction by leaving out those dimensions that show low variance.

## PRINCIPAL COMPONENTS ANALYSIS (2)

* Consider the sampled vectors $X^k$ ($k = 1, ..., N$) with $M = 0$ (if $M \neq 0$, construct a new set of vectors $Z^k = X^k - M$).

* We are looking for a unit vector $u$ on which the $X^k$ will be projected. The projection is:

$$p_k = X^k \cdot u = X^{k^T} u = u^T X^k$$

* The projection's mean is also zero:

$$\frac{1}{N} \sum_{k=1}^{N} p_k = \frac{1}{N} \sum_{k=1}^{N} u^T X^k = u^T \frac{1}{N} \sum_{k=1}^{N} X^k = 0.$$

* The projection's variance:

$$\sigma^2(u) = \frac{1}{N} \sum_{k=1}^{N} p_k^2 = \frac{1}{N} \sum_{k=1}^{N} \left( u^T X^k \right) \left[ X^{k^T} u \right] = u^T \Sigma u.$$

## PRINCIPAL COMPONENTS ANALYSIS (3)

* At the maximal point of variance:

$$\sigma^2(u + \Delta u) = \sigma^2(u)$$

$$\sigma^2(u + \Delta u) = (u + \Delta u)^T \Sigma (u + \Delta u)$$

* Ignoring second-order terms:

$$\sigma^2(u + \Delta u) = u^T \Sigma u + \Delta u^T \Sigma U + u^T \Sigma \Delta u$$

* Because $\Sigma$ is symmetric, $A^T \Sigma B = B^T \Sigma A$, and therefore:

$$\sigma^2(u + \Delta u) = u^T \Sigma u + 2\Delta u^T \Sigma u = \sigma^2(u) + 2\Delta u^T \Sigma u$$

* It can be concluded that:

$$\Delta u^T \Sigma u = 0$$

* The unit-vector constraint means:

$$(u + \Delta u)^T (u + \Delta u) = 1 \text{ or } \Delta u^T u = 0.$$

## PRINCIPAL COMPONENTS ANALYSIS (4)

* Using the technique of *Lagrange multipliers* one gets:

$$\Delta u^T \Sigma u - \lambda \Delta u^T u = 0 \text{ or } \Sigma u = \lambda u.$$

* So, the vector $u$ should be an eigenvector of the covariance matrix $\Sigma$. But:

$$\sigma^2(u) = u^T \Sigma u = u^T \lambda u = \lambda.$$

* So the largest variance is achieved when the eigenvector corresponding to the largest eigenvalue is chosen.

* Note that the eigenvectors of a symmetrical matrix are orthogonal and can be chosen as a base.

## PRINCIPAL COMPONENTS ANALYSIS (5)

* Take the eigenvalues in decreasing order ($\lambda_1 > \lambda_2 > ... > \lambda_n$) and construct a matrix $\Phi = \left[ u_1, u_2, ..., u_n \right]$ where each column is an eigenvector.

* One can then write: $\Sigma \Phi = \Phi \Lambda = \Lambda \Phi$ where $\Lambda$ is a diagonal matrix with the eigenvalues in decreasing order in the diagonal. This leads to: $\Phi^{-1} \Sigma \Phi = \Phi^T \Sigma \Phi = \Lambda$.

* So, after a transformation with $\Phi$, the covariance matrix becomes a diagonal matrix.

## DIMENSIONALITY REDUCTION (1)

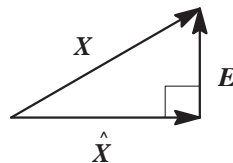* Expressing a vector $X$ as a linear combination of the eigenvectors gives: $X = \sum_{i=1}^{n} q_i u_i.$

* One can approximate $X$ by leaving out the last $n - m$ terms of the sum ($m < n$):

$$\hat{X} = \sum_{i=1}^{m} q_i u_i.$$

* The approximation error is then:

$$E = X - \hat{X} = \sum_{i=m+1}^{n} q_i u_i.$$

* Note that the error is always orthogonal to the approximating vector. This is the *principle of orthogonality.*

## DIMENSIONALITY REDUCTION (2)

* Consider the total variance of the approximating vector:

$$\sum_{i=1}^{m} \sigma_i^2 = \sum_{i=1}^{m} \lambda_i$$

where $\sigma_i$ is the variance in the $i$th dimension after projection on the base of eigenvectors $\Phi$.

* This means that it is indeed a good idea to order the eigenvalues from large to small and take as many vectors from $\Phi$ as desired.

## DIMENSIONALITY REDUCTION (3)

* Denote: $\Phi^m = \begin{bmatrix} \boldsymbol{u}_1, \dots, \boldsymbol{u}_m \end{bmatrix}$.

* $\hat{\boldsymbol{X}} = \sum_{i=1}^{m} q_i \boldsymbol{u}_i$, means that the $n$-dimensional $\boldsymbol{X}$ can now be represented by an $m$-dimensional vector $\boldsymbol{Q}$.

* So: $\hat{\boldsymbol{X}} = \Phi^m \begin{bmatrix} q_1 \\ \dots \\ q_m \end{bmatrix} = \Phi^m \boldsymbol{Q}$ .

* $\boldsymbol{Q}$ is found by: $Q = (\Phi^m)^T \boldsymbol{X}$ from the original vector $\boldsymbol{X}$.

For more information on PCA, consult:

[1] Haykin, S., Neural Networks, A Comprehensive Foundation, Prentice Hall International, Upper Saddle River, New Jersey, Second Edition, (1999).

[2] Jang, J.S.R., C.T. Sun and E. Mizutani, Neuro–Fuzzy and Soft Computing. A Computational Approach to Learning and Machine Intelligence, Prentice Hall, Upper Saddle River, NJ, (1997).

## HIGH-DIMENSIONAL SPACES

* The *sample covariance matrix* when the mean is zero, is defined as:

$$\Sigma = \frac{1}{N} \sum_{k=1}^{N} X^k (X^k)^T$$

* Create an $n \times N$ matrix $A$ composed of the sampled vectors:

$$A = \begin{bmatrix} X^1, X^2, \dots, X^N \end{bmatrix}$$

* Then the $n \times n$ covariance matrix can also be written as:

$$\Sigma = \frac{1}{N} A A^T.$$

* Finding the eigenvectors of $\Sigma$ for large $n$ is difficult, if not intractable. We consider here cases where $N < n$.

* Consider an eigenvector $\boldsymbol{v}$ of $A^T A$, an $N \times N$ matrix:

$$A^T A \boldsymbol{v} = \mu \boldsymbol{v}.$$

* Premultiply by $A$:

$$A A^T A \boldsymbol{v} = A \mu \boldsymbol{v} = \mu A \boldsymbol{v}.$$

* So, when $\boldsymbol{v}$ is an eigenvector of $A^T A$, $A\boldsymbol{v}$ is an eigenvector of $\Sigma$!

## EXAMPLE: EIGENFACES

* It considers images of faces represented by $256 \times 256 = 65536$ pixels.

* A direct application of the PCA to find the eigenvalues and eigenvectors of the covariance matrix would lead to a calculation involving $65536 \times 65536$ matrix!

* However, using the technique for high-dimensional spaces requires finding the eigenvectors for an $N \times N$ matrix, where $N$ is the number of samples, e.g. 16.

* The eigenvectors found form a base to represent faces, the so-called *face space*.

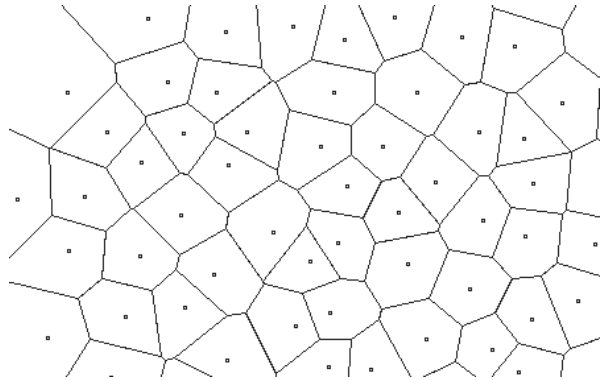* The coordinates in face space form a feature vector that can be used for face recognition.

## CLUSTERING: VECTOR QUANTIZATION

* *Clustering* is the process of representing a set of $N$ data points $X^i$ ($i = 1, \dots, N$) by a set of $M$ data points $Y^j$ ($j = 1, \dots, M, M < N$); the smaller set should be in some sense "representative" of the larger one.

* *Vector quantization* is an application of clustering in lossy data compression. The $M$ data points form a *codebook* that is available at the receiver side. Given a point $X^i$ to transmit, the sender finds the point $Y^j$ closest to $X^i$ from the codebook and simply sends the index $j$. The receiver reconstructs an approximation of the data by retrieving $Y^j$ from the codebook.

# VORONOI DIAGRAMS



* Subdivide the plane in *Voronoi* regions by perpendicular bisectors between neighboring pairs of points.

---

# $k$-MEANS CLUSTERING ALGORITHM (1)

* It finds $k$ mean vectors that represent a set of data points by means of $k$ classes.
* The algorithm is as follows:

  1) Start with cluster points $Y^j$ ($j = 1, ..., M$) and determine the corresponding Voronoi region $\mathcal{V}^j$.

  2) Determine the *centroids* of the data points $X^i$ ($i = 1, ..., N$) contained in each Voronoi region:

  $$\frac{1}{|\mathcal{V}^j|} \sum_{X^i \in V^j} X^i$$

  3) Assign the centroids to the cluster points $Y^j$ and repeat from Step 2 until convergence.

---

# $k$-MEANS CLUSTERING ALGORITHM (2)

Remarks:

* The algorithm will always converge, but not necessarily to the global minimum.
* Quality of final solution strongly depends on initial assignment of cluster points (e.g. random choice, or a choice based on principle components).

Source:

[3] Moon, T.K. and W.C. Stirling, "Mathematical Methods and Algorithms for Signal Processing", Prentice Hall, Upper Saddle River, (2000).