

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Slido event code: 1872043
- Homework 7 is due Friday
 - For Pr. 2, use dataset
Wine_data_v2.csv

Today's lecture

- d.o.f. and constraints in ANNs
 - Ex1: RBF networks
 - Ex2: More general ANNs

- Review of random vectors
 - Definitions
 - First and second order statistics
 - Multivariate Normal
 - Mahalanobis distance
 - Linear transformations
- Orthonormal transformation
- Whitening transformation
- Start Statistical Classification
 - Bayes Decision Theory (part 1)

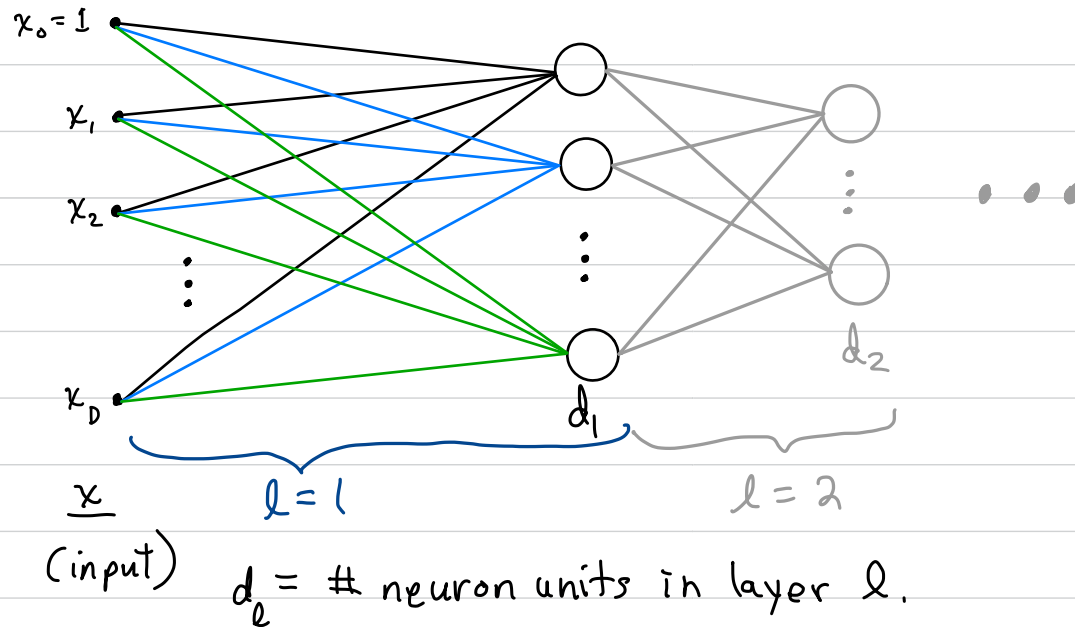
* p.2: activation fcn. notation changed to h post-lecture.

Complexity: degrees of freedom and constraints in ANN

Each weight that is varied during learning, is a d.o.f., unless it is "tied" or constrained in some way.

Other parameters that are optimized as a function of the data are also d.o.f.

Ex 1: 1-layer (fully connected) ANN:



* $l=1 \Rightarrow$ Assume activation functions h are given (not a fcn. of data)
 d.o.f. = $\# w_{jd}^{(1)}$ scalars = $(D+1) \cdot d_1$, if each $w_{jd}^{(1)}$ is not constrained or tied to other weights in some way.

- d.o.f. $\Rightarrow N_c = ?$ (1 layer ANN)

Consider a 1000×1000 color image as input. $D = ?$ $3 \cdot 10^6$

If $d_i = 200 \times 200$, with $l=1$ fully connected.

Then how many weights in $l=1$?

$$4 \cdot 10^4 \times 3 \cdot 10^6 \approx 10^{11}$$

- How many images would we need in our training set?

$$\Rightarrow N > (3-10) \text{ d.o.f.} \approx (3-10) \cdot 10^{11} \sim 10^{12} = 1 \text{ Trillion !}$$

Ex 2: ChatGPT

GPT 3; has 175 Billion parameters.
(Wikipedia has 4.2 billion words.)

GPT 4: 1.7 Trillion parameters

GPT 5: 2T - 5T

Will GPT run out of words to learn from?

Can reduce d.o.f. by:

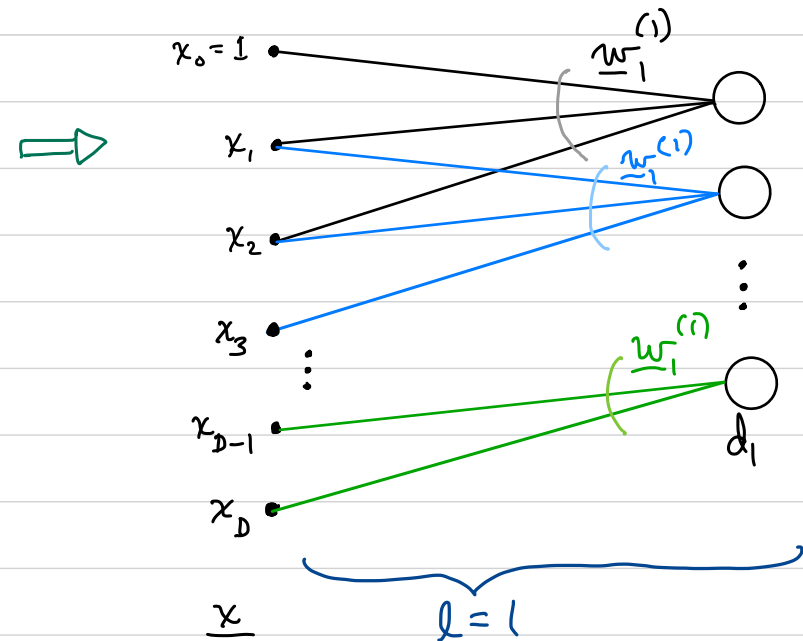
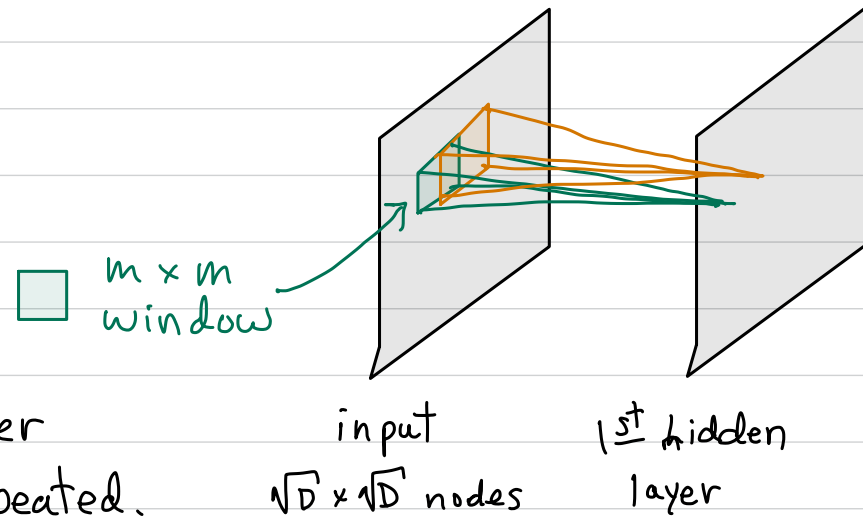
(1) Eliminate some interconnections
e.g., keep interconnections
"local".

(2) "Tie" some weights to each other
e.g., same set of weights is repeated.

(1) & (2) are used to reduce d.o.f.
in convolutional neural networks
(CNNs), commonly used for ML
on images and for computer vision.
[Bishop Sec.5.5.6; Goodfellow Ch.9]

(3) regularization can be used
to constrain the weights
in some other way.
[Bishop 5.5 describes some techniques.]

Example for images as input



Random Vectors and Their Properties

Ref: Bishop 2.3.0-2.3.3; Duda, Hart, and Stork A.4.8-A.5.2; 2.5

Definitions and Properties for 1 Random Vector

Let $\underline{x} = [x_1, x_2, \dots, x_D]^T$ in which the components of \underline{x} are random variables (r.v.).

The joint probability density function p is:

$$p(\underline{x}) = p(x_1, x_2, \dots, x_D)$$

in which p is a scalar. Note that:

$$\int p(\underline{x}) d\underline{x} \triangleq \iint \dots \int p(x_1, x_2, \dots, x_D) dx_1 dx_2 \dots dx_D = 1$$

Mean or Expected Value

$$\underline{m} = E\{\underline{x}\} = \int \underline{x} p(\underline{x}) d\underline{x}$$

$$m_i = \iint \dots \int x_i p(x_1, x_2, \dots, x_D) dx_1 dx_2 \dots dx_D = \int x_i \underbrace{\iint \dots \int}_{\text{all } i \neq k} p(x_i, \dots, x_D) dx_1 \dots dx_D$$

$$m_i = \int x_i p(x_i) dx_i$$

$$\text{in which } p(x_i) = \iint \dots \int p(x_1, x_2, \dots, x_D) dx_1 dx_2 \dots dx_{i-1} dx_{i+1} \dots dx_D$$

Autocorrelation and Cross-correlation Matrix

Autocorrelation matrix:

$$\begin{aligned}
 E\{\underline{x}\underline{x}^T\} &= E\left\{ \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \begin{bmatrix} x_1 & \cdots & x_D \end{bmatrix} \right\} \\
 &= E\left\{ \begin{bmatrix} x_1x_1 & x_1x_2 & x_1x_3 & \cdots & x_1x_D \\ x_2x_1 & & & & \\ x_3x_1 & & \ddots & & \vdots \\ \vdots & & & & \\ x_Dx_1 & \cdots & & x_Dx_D \end{bmatrix} \right\} \\
 &= \begin{bmatrix} E\{x_1x_1\} & E\{x_1x_2\} & E\{x_1x_3\} & \cdots & E\{x_1x_D\} \\ E\{x_2x_1\} & & & & \\ E\{x_3x_1\} & & \ddots & & \vdots \\ \vdots & & & & \\ E\{x_Dx_1\} & \cdots & & & E\{x_Dx_D\} \end{bmatrix}
 \end{aligned}$$

e.g.:
living area
#rooms

in which $E\{x_i x_j\} = \iint x_i x_j p(x_i, x_j) dx_i dx_j$.

Cross-correlation matrix is $E\{\underline{x}\underline{y}^T\}$.

Covariance Matrix

$$\begin{aligned}\underline{\underline{\Sigma}} &= E \left\{ (\underline{x} - \underline{m})(\underline{x} - \underline{m})^T \right\} \\ &= \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} & \cdots & \sigma_{1D} \\ \sigma_{21} & & & & \\ \sigma_{31} & & \ddots & & \vdots \\ \vdots & & & & \\ \sigma_{D1} & \cdots & & & \sigma_{DD} \end{bmatrix}\end{aligned}$$

in which $\sigma_{ij} = E \left\{ (x_i - m_i)(x_j - m_j) \right\}$, and $\sigma_{ii} = \sigma_i^2 = E \left\{ (x_i - m_i)^2 \right\}$.

Note that \uparrow covariance of x_i and x_j \uparrow variance of x_i (i th feature).

$$\begin{aligned}\underline{\underline{\Sigma}} &= E \left\{ (\underline{x} - \underline{m})(\underline{x} - \underline{m})^T \right\} \\ &= E \left\{ \underline{x} \underline{x}^T \right\} - E \left\{ \underline{m} \underline{x}^T \right\} - E \left\{ \underline{x} \underline{m}^T \right\} + E \left\{ \underline{m} \underline{m}^T \right\} \\ &= E \left\{ \underline{x} \underline{x}^T \right\} - \underline{m} E \left\{ \underline{x}^T \right\} - E \left\{ \underline{x} \right\} \underline{m}^T + \underline{m} \underline{m}^T \\ &= E \left\{ \underline{x} \underline{x}^T \right\} - \underline{m} \underline{m}^T - \cancel{\underline{m} \underline{m}^T} + \cancel{\underline{m} \underline{m}^T} \\ \underline{\underline{\Sigma}} &= E \left\{ \underline{x} \underline{x}^T \right\} - \underline{m} \underline{m}^T\end{aligned}$$

$\underline{\underline{\Sigma}}$ is symmetric and positive semi-definite.

We will assume $\underline{\underline{\Sigma}}$ is positive definite, so that $|\underline{\underline{\Sigma}}| > 0$.

$|\underline{\underline{\Sigma}}| = 0$ is a degenerate case; for example this will happen when:

$$\sigma_{ii} = 0 \text{ for some } i, \text{ or } x_i = \alpha x_j.$$

The covariance matrix is often normalized:

$$\underline{\underline{\Sigma}} = \underline{\underline{\Gamma}} \underline{\underline{\mathbf{R}}} \underline{\underline{\Gamma}}$$

in which

$$\underline{\underline{\Gamma}} = \begin{bmatrix} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ 0 & & & & \sigma_D \end{bmatrix}, \quad \underline{\underline{\mathbf{R}}} = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & r_{ji} & & \ddots & \\ & & & & 1 \end{bmatrix}$$

with

$$r_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$$

Note that $0 \leq |r_{ij}| \leq 1$. $\underline{\underline{\mathbf{R}}}$ is often called a “correlation matrix”, but this terminology can be confusing. Another term is “normalized covariance matrix”.

Definitions and Properties for 2 Random Vectors

Two random vectors \underline{x} and \underline{y} are defined as:

- Uncorrelated if $E\{\underline{x} \underline{y}^T\} = E\{\underline{x}\} E\{\underline{y}^T\}$
- Orthogonal if $E\{\underline{x}^T \underline{y}\} = 0$
- Independent if $p(\underline{x}, \underline{y}) = p(\underline{x}) p(\underline{y})$

As a result,

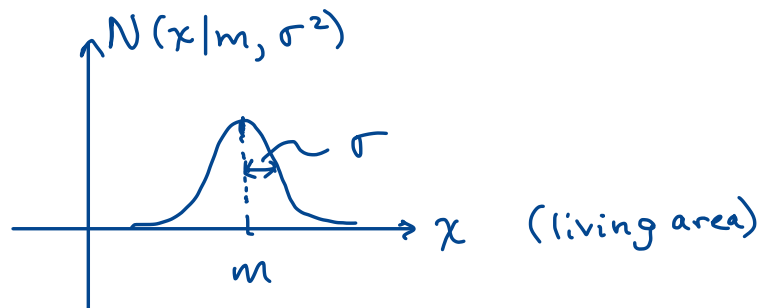
- (1) independent \Rightarrow uncorrelated
- (2) uncorrelated \nRightarrow independent, in general.
- (3) If $E\{\underline{x}\} = \underline{0}$ or $E\{\underline{y}\} = \underline{0}$, then uncorrelated \Rightarrow orthogonal.

Normal or Gaussian density [DHS 2.5]

Univariate case

$$p(x) = N(x|m, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-m}{\sigma}\right)^2\right\}$$

$$E\{x\} = m, \quad E\{(x-m)^2\} = \sigma^2$$



Multivariate case

$$p(\underline{x}) = N(\underline{x} | \underline{m}, \underline{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\underline{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} d_M^2(\underline{x}, \underline{m}, \underline{\Sigma}) \right\}$$

in which $|\underline{\Sigma}|$ = determinant of $\underline{\Sigma}$,

$$\begin{aligned} d_M^2(\underline{x}, \underline{m}, \underline{\Sigma}) &= (\underline{x} - \underline{m})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{m}) \\ &= \text{tr} \left\{ \underline{\Sigma}^{-1} (\underline{x} - \underline{m})(\underline{x} - \underline{m})^T \right\} \end{aligned}$$

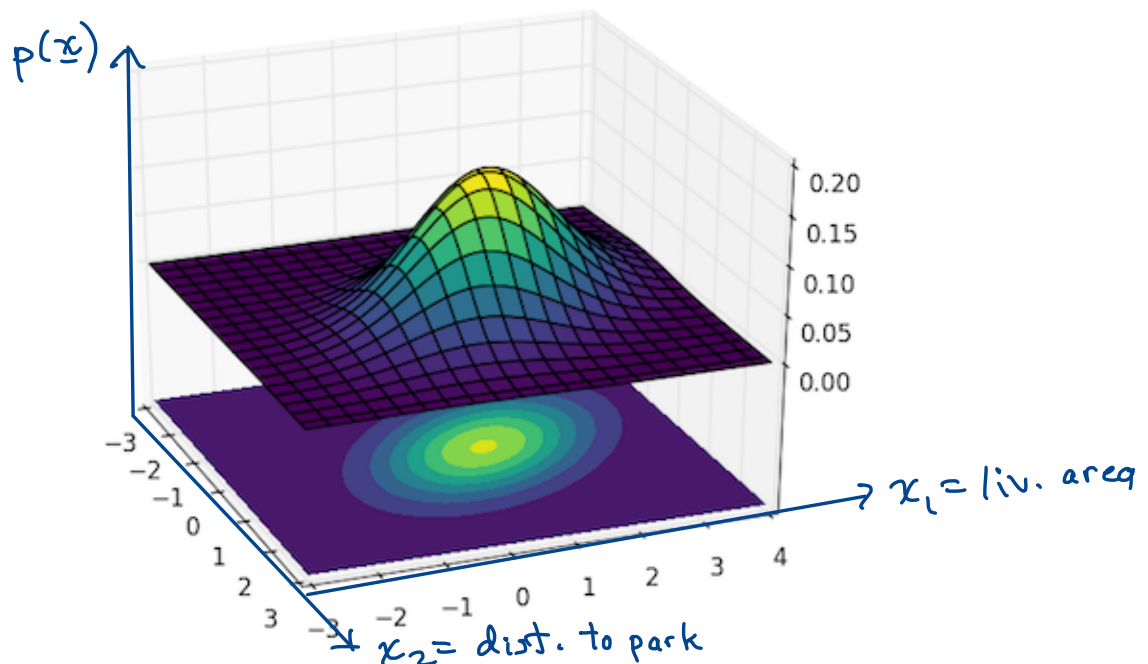
$\text{tr}\{\underline{A}\}$ = trace of \underline{A} , and

$d_M(\underline{x}, \underline{m}, \underline{\Sigma})$ = Mahalanobis distance between \underline{x} and \underline{m} .

Thus, the multivariate normal is a function of:

D mean values and

$\frac{D(D+1)}{2}$ variance values.



Plot of bivariate Normal [1]

[1] scipython.com/blog/visualizing-the-bivariate-gaussian-distribution/

Uncorrelated Features

If the x_i and x_j are uncorrelated $\forall j \neq i$, then for $i \neq j$:

$$\begin{aligned}\sigma_{ij} &= E\left\{(x_i - m_i)(x_j - m_j)\right\} \\ &= E\left\{x_i x_j\right\} - m_i m_j = E\left\{x_i\right\} E\left\{x_j\right\} - m_i m_j \\ &= 0\end{aligned}$$

Thus

$$\underline{\underline{\Sigma}} = \begin{bmatrix} \sigma_{11} & & & 0 \\ & \sigma_{22} & & \\ & & \ddots & \\ 0 & & & \sigma_{DD} \end{bmatrix}, \quad \underline{\underline{\Sigma}}^{-1} = \begin{bmatrix} \frac{1}{\sigma_{11}} & & & 0 \\ & \frac{1}{\sigma_{22}} & & \\ & & \ddots & \\ 0 & & & \frac{1}{\sigma_{DD}} \end{bmatrix}$$

and

$$\begin{aligned}d_M^2(\underline{x}, \underline{m}, \underline{\underline{\Sigma}}) &= (\underline{x} - \underline{m})^T \underline{\underline{\Sigma}}^{-1} (\underline{x} - \underline{m}) \\ \Rightarrow d_M^2(\underline{x}, \underline{m}, \underline{\underline{\Sigma}}) &= \sum_{i=1}^D \frac{(x_i - m_i)^2}{\sigma_{ii}}\end{aligned}$$

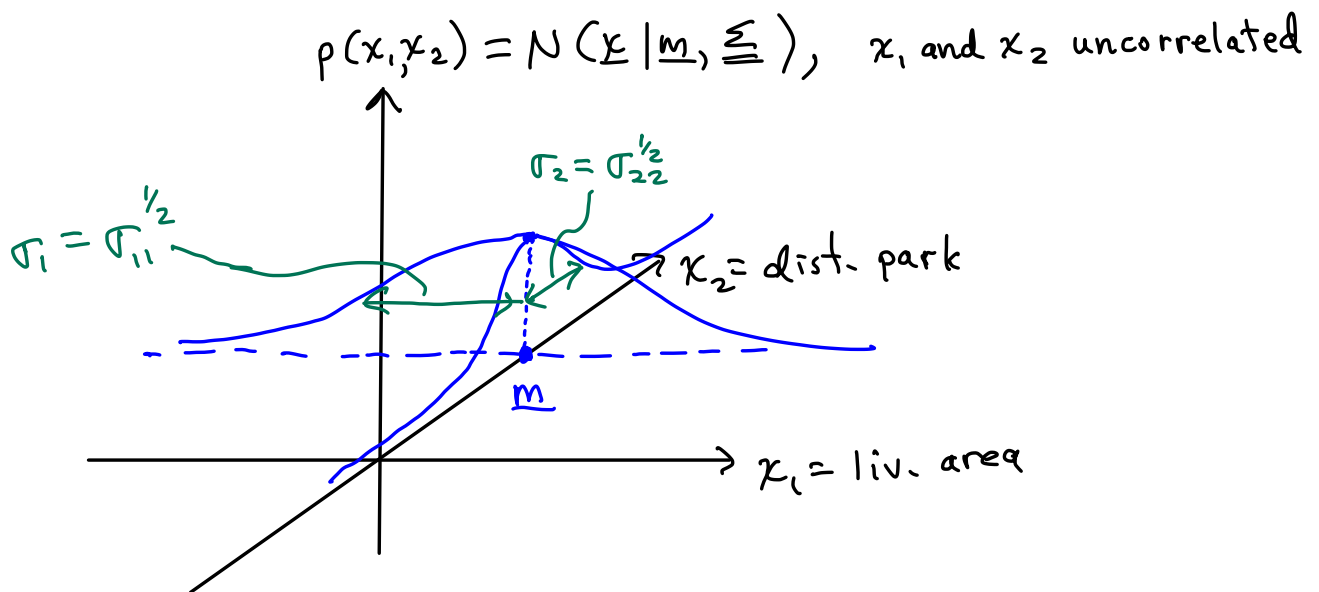
= Euclidean dist.² with auto-normalizing factor $\frac{1}{\sigma_{ii}} = \frac{1}{\sigma_i^2}$ for each feature x_i .

Also, for the normal case (with uncorrelated features):

$$\begin{aligned}
 p(\underline{x}) &= \frac{1}{(2\pi)^{D/2} \left(\prod_{i=1}^D \sigma_{ii} \right)^{1/2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^D \frac{(x_i - m_i)^2}{\sigma_{ii}} \right\} \\
 &= \prod_{i=1}^D \left[\frac{1}{(2\pi)^{1/2} \sigma_{ii}^{1/2}} \exp \left\{ -\frac{1}{2} \frac{(x_i - m_i)^2}{\sigma_{ii}} \right\} \right] \\
 &= \prod_{i=1}^D p(x_i) = \prod_{i=1}^D N(x_i | m_i, \sigma_{ii})
 \end{aligned}$$

$d_m^2(\underline{x}, \underline{m}, \underline{\Sigma})$

So uncorrelated \Rightarrow independent for the normal case.



Linear Transformation of Random Vectors

Let $\underline{y} = \underline{\underline{A}} \underline{x}$,

\underline{x} and \underline{y} are random vectors

\underline{x} in orig. feat. space
 \underline{y} in transformed
 feat. space.

Then $\underline{m}_y = \underline{\underline{A}} \underline{m}_x$,

$\underline{\underline{A}}$ is deterministic

And

$$\begin{aligned}\underline{\underline{\Sigma}}_y &= E \left\{ (\underline{y} - \underline{m}_y) (\underline{y} - \underline{m}_y)^T \right\} \\ &= E \left\{ (\underline{\underline{A}} \underline{x} - \underline{\underline{A}} \underline{m}_x) (\underline{\underline{A}} \underline{x} - \underline{\underline{A}} \underline{m}_x)^T \right\} \\ &= \underline{\underline{A}} E \left\{ (\underline{x} - \underline{m}_x) (\underline{x} - \underline{m}_x)^T \right\} \underline{\underline{A}}^T\end{aligned}$$

$$\Rightarrow \underline{\underline{\Sigma}}_y = \underline{\underline{A}} \underline{\underline{\Sigma}}_x \underline{\underline{A}}^T$$

Assuming $\underline{\underline{A}}$ is nonsingular,

$$\begin{aligned}d_M^2(\underline{y}, \underline{m}_y, \underline{\underline{\Sigma}}_y) &= (\underline{y} - \underline{m}_y)^T \underline{\underline{\Sigma}}_y^{-1} (\underline{y} - \underline{m}_y) \\ &= [\underline{\underline{A}} (\underline{x} - \underline{m}_x)]^T [\underline{\underline{A}} \underline{\underline{\Sigma}}_x \underline{\underline{A}}^T]^{-1} [\underline{\underline{A}} (\underline{x} - \underline{m}_x)] \\ &= (\underline{x} - \underline{m}_x)^T \underline{\underline{A}}^T (\underline{\underline{A}}^T)^{-1} \underline{\underline{\Sigma}}_x^{-1} \underline{\underline{A}}^{-1} \underline{\underline{A}} (\underline{x} - \underline{m}_x) \\ &= (\underline{x} - \underline{m}_x)^T \underline{\underline{\Sigma}}_x^{-1} (\underline{x} - \underline{m}_x) \\ &= d_M^2(\underline{x}, \underline{m}_x, \underline{\underline{\Sigma}}_x)\end{aligned}$$

\Rightarrow Mahalanobis distance is preserved under a linear transformation with nonsingular $\underline{\underline{A}}$.

Orthonormal Transformation

is a special case of a linear transformation.

Let $\underline{y} = \underline{\underline{E}}^T \underline{x}$, $\Rightarrow \underline{\underline{A}} = \underline{\underline{E}}^T$. Let $\underline{\underline{E}}^T \underline{\underline{E}} = \underline{\underline{I}} = \underline{\underline{E}} \underline{\underline{E}}^T$ ($\underline{\underline{E}}$ is orthonormal).

Also let $\underline{\underline{E}}$ = eigenmatrix of $\underline{\underline{\Sigma}}_x$: $\underline{\underline{E}} \triangleq \begin{bmatrix} \underline{e}_1 & \underline{e}_2 & \cdots & \underline{e}_D \end{bmatrix}$

in which \underline{e}_n is the n^{th} eigenvector of $\underline{\underline{\Sigma}}_x$, and $\{\underline{e}_n, n = 1, 2, \dots, D\}$ is orthonormal.¹

$$\underline{\underline{E}} \text{ satisfies: } \underline{\underline{\Sigma}}_x \underline{\underline{E}} = \underline{\underline{E}} \underline{\underline{\Lambda}}, \text{ in which } \underline{\underline{\Lambda}} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \cdots & 0 & \lambda_D \end{bmatrix}$$

and λ_n is the n^{th} eigenvalue of $\underline{\underline{\Sigma}}_x$.²

Note that: $\underline{\underline{\Sigma}}_y = \underline{\underline{A}} \underline{\underline{\Sigma}}_x \underline{\underline{A}}^T = \underline{\underline{E}}^T (\underline{\underline{\Sigma}}_x \underline{\underline{E}}) = \underline{\underline{E}}^T \underline{\underline{E}} \underline{\underline{\Lambda}} = \underline{\underline{\Lambda}}$. Thus $\underline{\underline{\Sigma}}_y = \underline{\underline{\Lambda}}$

and $\underline{\underline{E}}$ diagonalizes $\underline{\underline{\Sigma}}_x$.

$\underline{\underline{\Sigma}}_y = \text{diagonal} \Rightarrow$ the components of \underline{y} are uncorrelated.

Comments:

1. Orthonormal transformations are the basis of PCA.
2. It can be shown that Euclidean distance is preserved under an orthonormal transformation.

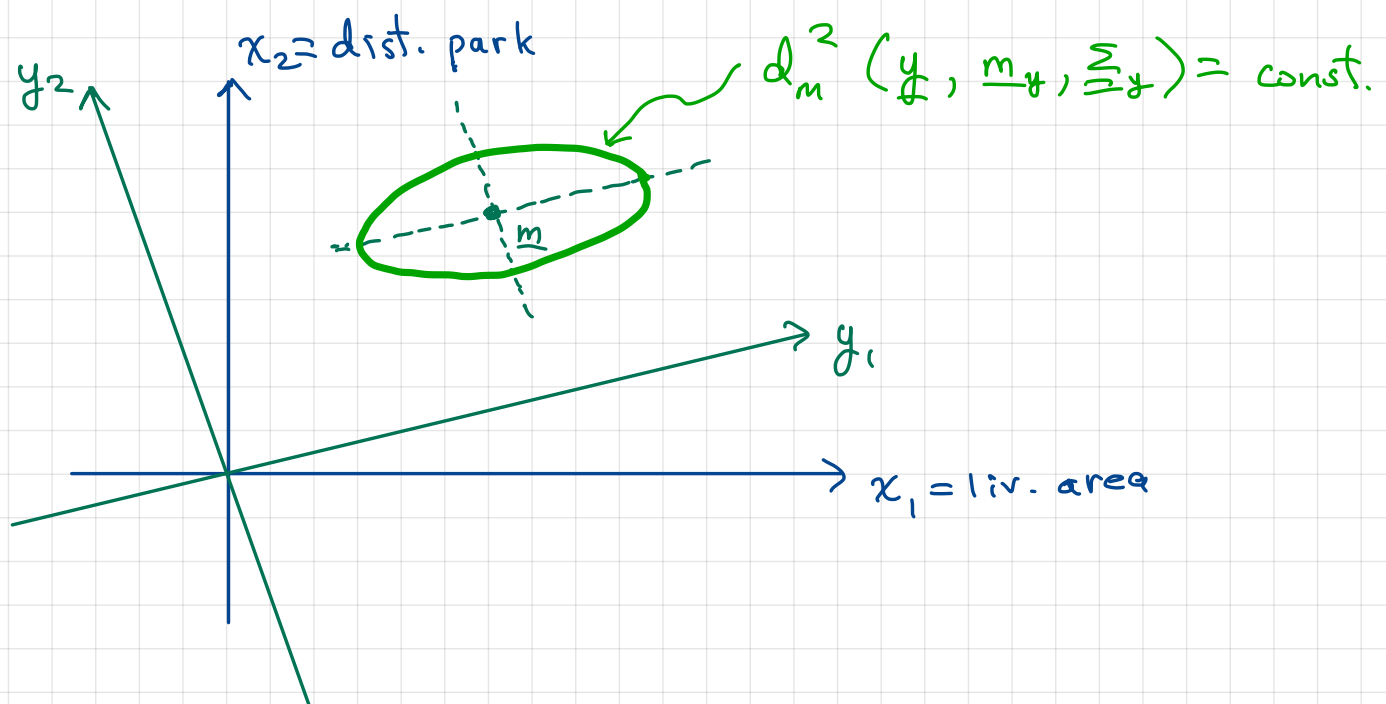
§

¹ Note that $\{\underline{e}_n, n = 1, 2, \dots, D\}$ is orthogonal, because $\underline{\underline{\Sigma}}_x$ is real and symmetric. The set can be made orthonormal by normalizing each eigenvector to unit length.

² All λ_n are real because $\underline{\underline{\Sigma}}_x$ is symmetric. All $\lambda_n > 0$ because we assumed $\underline{\underline{\Sigma}}_x$ is positive definite.

Ex:

p. 15



Whitening Transformation

Another special case of a linear transformation

$$\text{Let } \underline{y} = \underline{\Lambda}^{-1/2} \underline{E}^T \underline{x} \Rightarrow \underline{A} = \underline{\Lambda}^{-1/2} \underline{E}^T$$

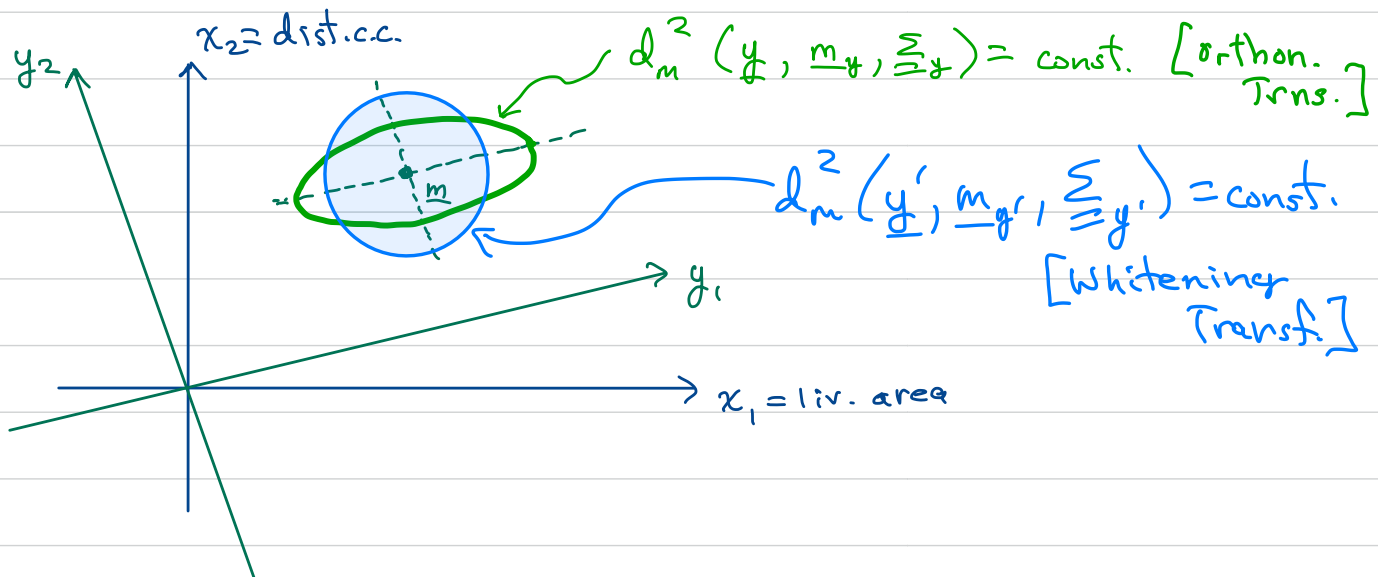
[in which \underline{E} = eigenmatrix of $\underline{\Sigma}_x$, and $\underline{\Lambda}$ = eigenvalue mtx. of $\underline{\Sigma}_x$

$$\underline{\Sigma}_y = \underline{A} \underline{\Sigma}_x \underline{A}^T = (\underline{\Lambda}^{-1/2} \underline{E}^T) \underline{\Sigma}_x (\underline{E} \underline{\Lambda}^{-1/2})$$

$$\underline{\Sigma}_y = \underline{\Lambda}^{-1/2} \underline{\Lambda} \underline{\Lambda}^{-1/2} = \underline{\Lambda}^{-1/2} \underline{\Lambda}^{1/2} = \underline{I} = \text{identity matrix}$$

$$\Rightarrow \sigma_{ii}^2 = 1, \sigma_{ij}^2 = 0 \text{ for } i \neq j.$$

\therefore The elements of \underline{y} are uncorrelated and have unit variance.



Comment: The whitening transformation is not orthonormal, and Euclidean distances are not preserved.

$$\|\underline{y}_1 - \underline{y}_2\|_2^2 \neq \|\underline{x}_1 - \underline{x}_2\|_2^2$$

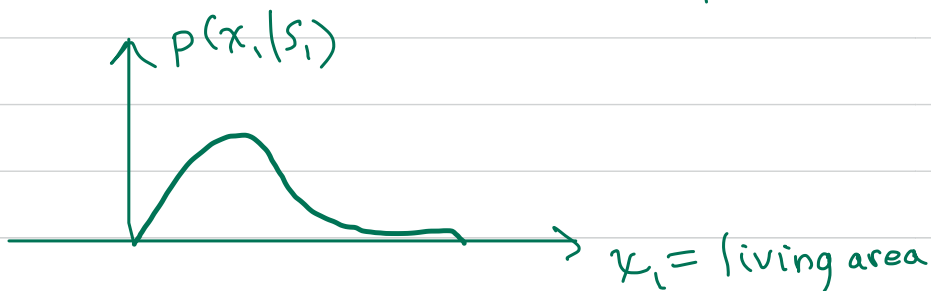
STATISTICAL CLASSIFICATION

Assume: data pts. \underline{x}_i (and unknowns \underline{x}) are drawn i.i.d. from $p(\underline{x})$.

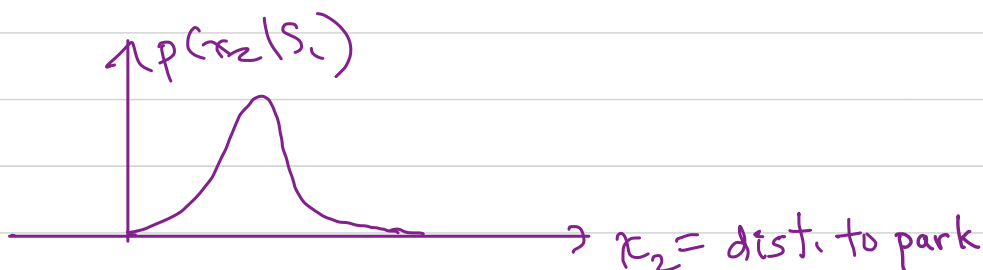
Bayes Decision Theory [Bishop 1.5]

Consider $p(\underline{x} | S_i) =$ class-conditional density
e.g.:

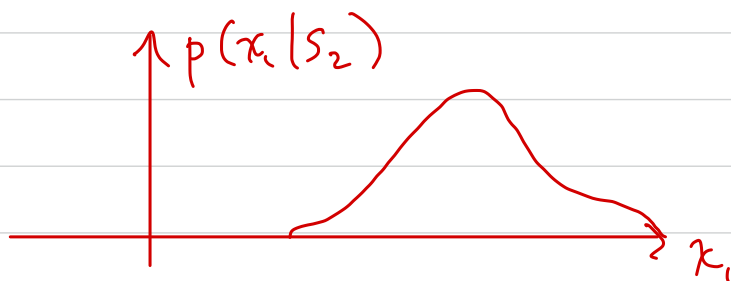
$$p(x_1 | S_1) = p(\text{living area} | \text{price increase})$$



$$p(x_2 | S_1) = p(\text{distance to park} | \text{price increase})$$



$$p(x_1 | S_2) = p(\text{living area} | \text{price decrease})$$



Assume: $p(\underline{x} | S_i)$ is known $\forall i$

Priors: $P(S_i)$

e.g.: $P(S_1) = P(\text{price increase})$, without knowledge of \underline{x} .

Assume: $P(S_i)$ is known $\forall i$

If $P(S_i)$ are unknown, can estimate by:

$$\hat{P}(S_i) = \frac{N_i}{N}, \text{ in which } N = \sum_{i=1}^C N_i$$

Bayes minimum-error classifier ($C=2$ classes)

Goal: choose decision boundary and regions that minimize $P(\text{error}) = P_e$

Let "e" denote "error"

$$P_e = P(e, S_1) + P(e, S_2)$$

$$= \underbrace{P(e | S_1)}_{\int_{\Gamma_2} p(\underline{x} | S_1) d\underline{x}} P(S_1) + \underbrace{P(e | S_2)}_{\int_{\Gamma_1} p(\underline{x} | S_2) d\underline{x}} P(S_2)$$

$$(1) \therefore P_e = P(S_1) \int_{\Gamma_2} p(\underline{x} | S_1) d\underline{x} + P(S_2) \int_{\Gamma_1} p(\underline{x} | S_2) d\underline{x}$$