

Multivariate Data

 $\underline{X}$  -  $d$ -columns -  $d$  variables

(inputs / features / attributes)

 $N$ -rows - iid.

observations / examples / instances

$$\underline{X} = \begin{bmatrix} X_1^1 & X_2^1 & \dots & X_d^1 \\ X_1^2 & & \dots & X_d^2 \\ \vdots & & & \vdots \\ X_1^N & \dots & \dots & X_d^N \end{bmatrix}$$

EG. Loan application -

customer observation vector - age, marital status, income, etc.

Typically, these variables are correlated

if not, great simplification.

Goals: simplification / summarization - explain data by few parameters.  
exploratory, - generate hypothesis about data.

Parameter Estimation

$$E[\underline{x}] = \underline{\mu} = [\mu_1 \dots \mu_d]$$

Variance of  $X_i$  is  $\sigma_i^2$ , covariance  $\bar{\sigma}_{ij} = \text{Cov}(X_i, X_j)$   
 $= E[(X_i - \mu_i)(X_j - \mu_j)^T]$ . $\Sigma$  is a matrix,  $\Sigma_{ii} = \sigma_i^2$ ,  $\Sigma_{ij} = \bar{\sigma}_{ij}$ .

$$\Sigma = E[(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^T] = E[\underline{x}\underline{x}^T] - \underline{\mu}\underline{\mu}^T$$

$$\text{Correlation} \quad \text{Corr}(X_i, X_j) = \rho_{ij} = \frac{\bar{\sigma}_{ij}}{\sqrt{\sigma_i^2 \sigma_j^2}}$$

If variables are indep., then correlation and covariance is zero

(2)

## Parameters Estimation (cont)

Given multivariate sample  
Sample mean.  $\bar{m} = \frac{\sum_{t=1}^N \bar{x}_t}{N}$ ,  $m_i = \frac{\sum_{t=1}^N \bar{x}_{it}}{N}$

Sample covariance  $\hat{\Sigma}$  with entries:

$$s_i^2 = \frac{1}{N-1} \sum_{t=1}^N (\bar{x}_{it} - m_i)^2 / N$$

$$s_{ij} = \frac{1}{N-1} \sum_{t=1}^N (\bar{x}_{it} - m_i)(\bar{x}_{jt} - m_j) / N.$$

Sample correlation.  $r_{ij} = s_{ij} / s_i s_j$ .

Missing Values — some variables may be missing in observations.

imputation — estimate the missing variable

— mean imputation — substitute the mean of other observations.

— imputation by regression — predict these values from existing ones.

but — data may not be missing at random

if so, need to model why it is missing  
(e.g. heights of Americans — Army data)

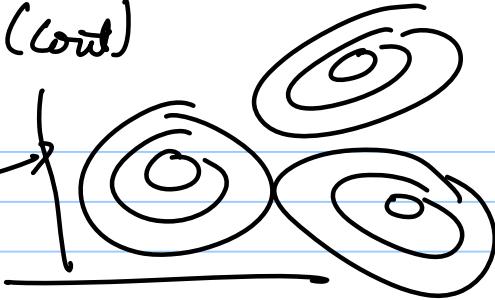
### Multivariate Normal

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Mahalanobis distance  $(x-\mu)^T \Sigma^{-1} (x-\mu)$

Level sets  $(x-\mu)^T \Sigma^{-1} (x-\mu) = c^2$  are ellipsoids centered on  $\mu$ , shape & orientation determined by  $\Sigma$ .

(5)

Multivariate Normal (cont)Bivariate case ( $d=2$ )  
Level sets.If  $\underline{x} \sim N_d(\mu, \Sigma)$ then each dimension of  $\underline{x}$  is a univariate normal (converse not necessarily true).Special case - if the components of  $\underline{x}$  are independent then the correlations (covariates) are 0.

$$P(\underline{x}) = \prod_{i=1}^d P_i(x_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} \sum_{i=1}^d \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2}.$$

(Note: we can always pick a coordinate rep for which this is true  $\rightarrow$  by rotating space to diagonal  $\Sigma$   
 $\Sigma$  e.g.  $\Sigma \rightarrow \varPhi^T \Sigma \varPhi$   $\varPhi$  rotates  
 $\varPhi$  can be obtained from eigenvectors/eigenvalues of  $\Sigma$ )

Important Propertythe projection of a  $d$ -dim normal onto a vector  $\omega$  is also a bivariate. $s_\omega$  is the projection onto a subspace  $\underline{\omega}$ .

$$\underline{\omega}^T \underline{x} = \omega_1 x_1 + \dots + \omega_d x_d \sim N(\omega^T \mu, \omega^T \Sigma \omega)$$

$$\underline{\omega}^T \underline{x} \sim N_k(\underline{\omega}^T \underline{\mu}, \underline{\omega}^T \underline{\Sigma} \underline{\omega}) \quad \underline{\omega} \text{ d x k matrix. rank } k < d.$$

(4)

## Multivariate Classifier

$$P(\underline{x} | C_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i)}$$

(analytic approach)

### Discriminant function

$$g_i(\underline{x}) = \log P(\underline{x} | C_i) + \log P(C_i)$$

$$g_i(\underline{x}) = -\frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \log P(C_i)$$

Given training samples  $\mathcal{X} = [\underline{x}^t, r^t]$

Estimate the distribution:  $\hat{P}(C_i) = \bar{r}_i^t / N$

$$\underline{\mu}_i = \bar{\underline{r}}_i^t \underline{x}^t / \bar{r}_i^t$$

$$\underline{\Sigma}_i = \bar{\underline{r}}_i^t (\underline{x}^t - \underline{\mu}_i)(\underline{x}^t - \underline{\mu}_i)^T / \bar{r}_i^t$$

Plug into discriminant function:

$$g_i(\underline{x}) = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \log \hat{P}(C_i)$$

### Quadratic Discriminant (dropping constant term)

Can be written as

$$g_i(\underline{x}) = \underline{x}^T \underline{w}_i \underline{x} + \underline{w}_i^T \underline{x} + w_{i0}$$

$$\underline{w}_i = -\frac{1}{2} \Sigma_i^{-1}, \quad \underline{w}_i = \Sigma_i^{-1} \underline{\mu}_i,$$

$$w_{i0} = -\frac{1}{2} \underline{\mu}_i^T \Sigma_i^{-1} \underline{\mu}_i - \frac{1}{2} \log |\Sigma_i| + \log \hat{P}(C_i).$$

No. parameters:  $K \cdot d$  mean +  $K \cdot d(d+1)/2$  covariance

(5)

## Multivariate Classification (cont.)

For large  $d$ , samples are small,  $\underline{\Sigma}_i$  may be singular and  $\underline{\Sigma}^{-1}$  won't exist.  
 May need to decrease dimensionality (next lecture)  
 Or pool data and estimate common covariance matrix:  $\underline{\Sigma} = \sum_i \hat{P}(C_i) \underline{\Sigma}_i$ .

$$\text{discriminant reduced to } g_i(\underline{x}) = -\frac{1}{2} (\underline{x} - \underline{m}_i)^T \underline{\Sigma}^{-1} (\underline{x} - \underline{m}_i) + \log \hat{P}(C_i)$$

Further simplification

assume off-diagonal elements covariance are zero.

naiive Bayes classifier:  $P(x_j | C_i)$  univariate Gaussian

$$g_i(\underline{x}) = -\frac{1}{2} \sum_{j=1}^d \left( \frac{x_j - m_{ij}}{s_j} \right)^2 + \log \hat{P}(C_i)$$

$s_j$  is standard deviation of  $j^{\text{th}}$  component  
 $m_{ij}$  is mean of  $j^{\text{th}}$  component of  $i^{\text{th}}$  class.

$\rightarrow$  reduces complexity of  $\underline{\Sigma}$  from  $O(n^2)$  to  $O(d)$ .

Extreme simplification: assume all variances are the same. Then Mahalanobis distance = Euclidean dist.

$$g_i(\underline{x}) = -\frac{1}{2} \sum_{j=1}^d (x_j - m_{ij})^2 + \log \hat{P}(C_i).$$

$$\text{If priors are equal} - g_i(\underline{x}) = -\|\underline{x} - \underline{m}_i\|^2$$

nearest mean classifier.  $\rightarrow$  prototype/template

$$g_i(\underline{x}) = \underline{w}_i^T \underline{x} + w_{i0}, \quad \underline{w}_i = \underline{m}_i, \quad w_{i0} = -\frac{1}{2} \|\underline{m}_i\|^2$$

'6)

## Tuning Complexity

As before, for polynomial regression, we can choose between models with different complexity.

→ Full Coverage  $d + \frac{1}{2}d(d+1)$  parameters

→ Indep. Sane Variables  $d + 1$

Same procedure as before:

(e.g. Cross-validation, regularization, etc.)

## Discrete Features:

Discrete attributes - color  $\in \{\text{red, blue, green, ...}\}$ ,  
pixel  $\in \{0, 1\}$

Suppose  $x_j$  are binary (Bernoulli)

$$p_{Cj} = p(x_j=1 | C_i)$$

$$P(x | C_i) = \prod_{j=1}^d p_{Cj}^{x_j} (1-p_{Cj})^{(1-x_j)}$$

$$\text{Discriminant } g_i(x) = \log P(x | C_i) + \log P(C_i)$$

$$= \sum_j (x_j \log p_{ij} + (1-x_j) \log (1-p_{ij})) + \log P(C_i)$$

$$\text{Estimator } \hat{p}_{ij} = \sum_t x_j^t r_i^t \underset{\text{linear in } x}{\underline{}}$$

General Case: multinomial  $x_j$  chosen from  $(v_1, \dots, v_n)$

$$p_{ijk} = \text{prob}(X_j \text{ from class } C_i \text{ takes value } v_k)$$

$$p_{ijk} = p(z_{jk}=1 | C_i) = p(x_j=v_k | C_i)$$

(7)

Discrete Features (cont)

$$z_{jk}^t = \begin{cases} 1, & \text{if } x_j^t = k \\ 0, & \text{otherwise} \end{cases}$$

If attributes are independent

$$P(x|C_i) = \prod_{j=1}^J \prod_{k=1}^{n_j} P_{ijk}$$

discretization  $g_i(x) = \sum_j \sum_k z_{jk} \log p_{ijk} + \log \hat{p}(C_i)$

Max like estimate for  $p_{ijk}$  is  $\hat{p}_{ijk} = \sum_t z_{jk}^t r_i^t / \sum_t r_i^t$

Multivariate Regression

$$r^t = g(\underline{x}^t | w_0, \dots, w_d) + \epsilon = w_0 + w_1 x_1^t + \dots + w_d x_d^t + \epsilon.$$

least squares:

$$E(w_0, \dots, w_d | \underline{x}) = \frac{1}{2} \sum_t (r^t - w_0 - w_1 x_1^t - \dots - w_d x_d^t)^2.$$

$$\underline{\underline{X}} = \begin{bmatrix} 1 & x_1^1 & \dots & x_d^1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^n & \dots & x_d^n \end{bmatrix}, \underline{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}, \underline{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix}$$

minimize wrt. w's

$$\underline{\underline{X}}^T \underline{\underline{X}} \underline{w} = \underline{\underline{X}}^T \underline{r}.$$

$$\hat{\underline{w}} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{r}.$$