

GAUSSIAN MODELS

MVN: $N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$

• EXPONENT: MAHALANOBIS DISTANCE BETWEEN DATA x AND MEAN VECTOR μ

• EIGEN DECOMPOSITION: $\Sigma = U \Lambda U^T$

EIGEN VECTORS ARE 'ELLIPSE' DIRECTION

MLE: $\hat{\mu} = \bar{x}$ $\hat{\Sigma} = \frac{1}{N} \left(\sum x_i x_i^T \right) - \bar{x} \bar{x}^T$ (EMPIRICAL MEAN AND COVARIANCE)

- HAS MAX ENTROPY AMONG DISTRIBUTIONS WITH SPECIFIED COVARIANCE Σ

GAUSSIAN DISCRIMINANT ANALYSIS

DEFINES CLASS CONDITIONAL DENSITIES IN A GENERATIVE CLASSIFIER $P(x|y=c, \theta) = N(x|\mu_c, \Sigma_c)$ • ALL CDD ARE GAUSSIAN

• Σ DIAGONAL \rightarrow GAUSSIAN NAIVE BAYES (FEATURES ARE c_i) • $\hat{y}(x) = \text{ARGMAX}_c [\log P(y=c|\pi) + \log P(x|\theta_c)]$

- MEASURES DISTANCE OF x FROM CENTER OF EACH CLASS μ_c , MAHALANOBIS DISTANCE \rightarrow 'NEAREST CENTROID CLASSIFIER'
- NORMALLY QUADRATIC, LINEAR WHEN Σ_c ARE TIED OR SHARED, $\Sigma_c = \Sigma$ • MIXED MODELS HAVE SAME FORM AS LOGISTIC REGRESSION

• REGULARIZED LDA $\Sigma = \lambda \text{DIAG}(\Sigma_{MVE}) + (1-\lambda) \Sigma_{MVE}$ • DIAGONAL LDA $\Sigma_c = \Sigma$, THEN DIAGONAL COV. MATRIX = TO PCA, $\lambda=1$ POOL ACROSS CLASSES WORKS BETTER IN HIGH DIMENSIONS

• NEAREST SHANNON CENTROIDS: MAP FOR DIAGONAL LDA WITH SPARCITY INDUCING PRIOR (LAPLACE)

CLASS SPECIFIC FEATURE MEAN $\mu_{c,j} = \mu_j + \Delta_{c,j}$ • RELEVANT TO LOGISTIC REGRESSION

MVN MARGINALS:

$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$, $\Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}$ $\begin{cases} P(x_1) = N(x_1|\mu_1, \Sigma_{11}) \\ P(x_2) = N(x_2|\mu_2, \Sigma_{22}) \end{cases}$

MVN

CONDITIONALS: $P(x_1|x_2) = N(x_1|\mu_{1|2}, \Sigma_{1|2})$ $\mu_{1|2} = \Sigma_{112}(\Lambda_{11}/\mu_1 - \Lambda_{12}(x_2 - \mu_2))$

$\Sigma_{1|2} = \Lambda_{11}^{-1}$

2D MVN $P(x_1|x_2) = N \left(x_1 | \mu_1 + \frac{\rho \sigma_1 \sigma_2}{\sigma_2^2} (x_2 - \mu_2), \sigma_1^2 - \frac{(\rho \sigma_1 \sigma_2)^2}{\sigma_2^2} \right)$

GAUSSIAN INTERPOLATION:

$x_j = \frac{1}{2} (x_{j-1} + x_{j+1}) + \epsilon_j$, $\epsilon_j \sim N(0, 1/\lambda \cdot I)$

D SUBINTERVALS

$Lx = \epsilon$: L IS FINITE DIFFERENCES MATRIX

- USED FOR DATA IMPUTATION

Information Form of MVN:

$$\Lambda = \Sigma^{-1}$$

$$\xi = \Sigma^{-1} \mu$$

$$N_c(x|\xi, \Lambda) = (2\pi)^{-D/2} |\Lambda|^{1/2} \exp \left[-\frac{1}{2} \left(x^T \Lambda x + \xi^T \Lambda^{-1} \xi - 2x^T \xi \right) \right]$$

- Λ, ξ canonical parameters
- μ, σ moment parameters

- CONDITIONING IS EASIER! $P(x_1|x_2) = N_c(x_1|\xi_1 - \Lambda_{12}x_2\Lambda_{11}^{-1})$

- MULTIPLYING IS EASIER! $N_c(\xi_1, \Lambda_1) \cdot N_c(\xi_2, \Lambda_2) = N_c(\xi_1 + \xi_2, \Lambda_1 + \Lambda_2)$

- MARGINALIZATION EASIER IN MOMENT FORM

LINEAR GAUSSIAN SYSTEMS

PRIOR

$$P(x) = N(x|\mu_x, \Sigma_x)$$

POSTERIOR

$$P(y|x) = N(y|Ax+b, \Sigma_y)$$

LIKELIHOOD

$$P(x|y) = N(x|\mu_{x|y}, \Sigma_{x|y})$$

$$\Sigma_{x|y}^{-1} = \Sigma_x^{-1} + A^T \Sigma_y^{-1} A, \quad \mu_{x|y} = \Sigma_{x|y} \left[A^T \Sigma_y^{-1} (y-b) + \Sigma_x^{-1} \mu_x \right]$$

DATA / NORMALIZATION

$$P(y) = N(y|A\mu_x+b, \Sigma_y + A\Sigma_x A^T)$$

- USEFUL TO INFER SCALARS AND VECTORS FROM NOISY DATA, INTERPOLATE DATA

• TIM HONOR RECOGNITION

• WISHART DISTRIBUTION

WISHART DISTRIBUTION

USED TO MODEL UNCERTAINTY IN COVARIANCE MATRICES Σ OR Λ , GENERALIZES Γ TO POSITIVE DEFINITE MATRICES, OR MULTIM OM OF χ^2

$$W_1(\Lambda|S, \nu) = \frac{1}{Z_{W_1}} |\Lambda|^{(\nu-D-1)/2} \exp \left(-\frac{1}{2} \text{TR}(\Lambda S^{-1}) \right)$$

ν = DOF

• ARE CONJUGATE PRIOR OF Λ

S = SCALE MATRIX

• W_1 IS DISTRIB OF SAMPLE COV. MATRIX OF MVN

$$Z_{W_1} = \text{NORMALIZATION} = 2^{\nu D/2} \Gamma_D(\nu/2) |S|^{\nu/2}$$

$$\Gamma_D(a) = \text{MULTIVARIATE GAMMA}$$

- ARE DISTRIBUTION OVER MATRICES

- $X_i \sim N(0, \Sigma) \rightarrow S = \sum_{i=1}^{\nu} X_i X_i^T$ IS WISHART $S \sim W_1(\Sigma, \nu)$, $E[S] = \nu \Sigma$, MEAN = νS

- $D=1 \rightarrow W_1(\lambda|s^{-1}, \nu) = \text{GA}(\lambda|\frac{\nu}{2}, \frac{s}{2})$

INVERSE WISHART

$$IW(\xi|S, \nu) = \frac{1}{Z_{IW}} |\xi|^{-(\nu+D+1)} \exp \left(-\frac{1}{2} \text{TR}(S^{-1} \xi^{-1}) \right) \quad Z_{IW} = |S|^{-\nu/2} 2^{\nu D/2} \Gamma_D(\nu/2)$$

- $D=1 \quad IW(\sigma^2|s^{-1}, \nu) = IG(\sigma^2|\nu/2, s/2)$

INVERSE WISHART \rightarrow INVERSE GAMMA

• CONJ PRIOR OF Σ OF MVN

INFERRING MVN PARAMETERS

$$x_i \sim N(\mu, \Sigma)$$

• POSTERIOR FOR μ

LIKELIHOOD $P(D|\mu) = N(\bar{x} | \mu, \frac{1}{N} \Sigma)$

GAUSSIAN PRIOR $P(\mu) = N(\mu | \mu_0, V_0)$

POSTERIOR (ALSO GAUSSIAN) $P(\mu | D, \Sigma) = N(\mu | m_N, V_N)$

IF UNINFORMATIVE PRIOR \rightarrow POSTERIOR MEAN \approx MLE

$$V_N^{-1} = V_0^{-1} + N \Sigma^{-1}$$

$$m_N = V_N (\Sigma^{-1} (N \bar{x}) + V_0^{-1} \mu_0)$$

• POSTERIOR FOR Σ

as d

$$P(\Sigma | D, \mu) = IW(\Sigma | S_N, V_N)$$

$$V_N = V_0 + N$$

$$S_N^{-1} = S_0 + S/N \text{ SCALAR MATRIX}$$

$$\hat{\Sigma}_{MAP} = \frac{S_0 + S_N}{N_0 + N}$$

REWRITE AS $\hat{\Sigma}_{MAP} = \lambda \hat{\Sigma}_0 + (1-\lambda) \hat{\Sigma}_{MLE}$ FOR SHRINKAGE / REGULARIZED ESTIM.

BAYESIAN T-TEST

$$P(\mu | D) = T\left(\mu | \bar{x}, \frac{s^2}{N}, N-1\right), \quad T = \frac{\bar{x} - \mu_0}{s/\sqrt{N}}, \quad P(\mu | D) = 1 - F_{N-1}(t)$$

- HERE μ IS UNKNOWN AND \bar{x} IS FIXED, CONTRARY TO FREQUENTIST FRAMEWORK

- SAME FORM AS FREQUENTIST T-TEST IF UNINFORMATIVE PRIOR

FISHER'S LDA

PROJECT 2 IN LOWER DIM. SPACE $Z = WX$, $W = L \times D$

↑

USUAL DA IS PROBLEMATIC IN HIGH DIMENSIONS. REDUCE DIMENSIONALITY. IE PCA, BUT SINCE IT'S UNSUPERVISED ITS RESULTS ARE NOT NECESSARILY OPTIMAL FOR CLASSIFICATION.

• **IDEA:** FIND W SO TO HAVE Z BE OPTIMALLY CLASSIFIED WITH A GAUSSIAN CLASS CONDITIONAL DENSITY. GAUSSIANTY IS OK SINCE IT'S LINEAR COMBS OF STUFF

→ BECAUSE RANK OF Σ_B IS $C-1$

• **LDA** REDUCES TO AT MOST $L \leq C-1$ DIMENSIONS; IN TWO-CLASS CASE W IS A SINGLE VECTOR

• **2-CLASS CASE** $\mu_1 = \frac{1}{N_1} \sum x_1$, $\mu_2 = \frac{1}{N_2} \sum x_2$ $m_H = W^T \mu_H$ PROJECTION OF MEAN ON LINE W . FIND W TO MAXIMIZE DISTANCE BETWEEN MEANS
 *** $S_B = \lambda S_W W \rightarrow$ GENERALIZED EIGENVALUE PROBLEM $\xrightarrow[SHORCUT]{2 \text{ CLASS}}$ $W = S_W^{-1} (\mu_2 - \mu_1)$ $J(W) = \frac{(\mu_2 - \mu_1)^2}{S_1^2 + S_2^2}$
 $S_B =$ BETWEEN CLASS SCATTER, $S_W =$ WITHIN-CLASS SCATTER

GENERAL: FIND EIGEN VECTORS OF $S_W^{-1} S_B$

• **HIGHER DIMENSIONS, MORE CLASSES** FIND MATRIX W TO MAXIMIZE $J(W) = \frac{W \Sigma_B W^T}{W \Sigma_W W^T}$ $W = \sum_{i=1}^{C-1} \lambda_i^{-1/2} U_i$ U IS L TOP EIGENVECTORS OF $\sum_{i=1}^{C-1} \lambda_i^{-1/2} \Sigma_B \lambda_i^{-1/2}$
 $\Sigma_{B,W} =$ BETWEEN / WITHIN CLASS COV. MATRICES

EXTENSIONS:

- **HETEROSKEDASTIC LDA:** Σ_C ARE NOT DIAGONAL AND EQUAL = FLOA
- **MULTIPLE LDA:** EACH CLASS HAS OWN PROJECTION MATRIX

$$S_B = \sum_i (\mu_i - \mu)(\mu_i - \mu)^T$$

$$S_W = \sum_i \sum_j p_i(x_j | \mu_i) (x_j - \mu_i)^T = \text{AUTOV} \times P_{\text{CLASS}}$$

- **ALGO:** ~~CONST~~ FIND FULL COV
 SEVERAL CLASS DATA
 FOR EACH CLASS
 FIND CLASS COV
 SUM TO GET S_W

$$S_B = C - S_W$$

FIND EIGEN

PICK TOP EIGEN $\rightarrow W_S$

PROJECT ORIGINAL DATA WITH W