

CS 340 Lec. 18: Multivariate Gaussian Distributions and Linear Discriminant Analysis

AD

March 2011

Multivariate Gaussian

- Consider data $\{\mathbf{x}^i\}_{i=1}^N$ where $\mathbf{x}^i \in \mathbb{R}^D$ and we assume they are independent and identically distributed.
- A standard pdf used to model multivariate real data is the multivariate Gaussian or normal

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} \underbrace{(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}_{\text{Mahalanobis distance}}\right). \end{aligned}$$

- It can be shown that $\boldsymbol{\mu}$ is the mean and $\boldsymbol{\Sigma}$ is the covariance of $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$; i.e.

$$\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu} \text{ and } \text{cov}(\mathbf{X}) = \boldsymbol{\Sigma}.$$

- It will be used extensively in our discussion on unsupervised learning and can also be used for generative classifiers (i.e. discriminant analysis).

- When $D = 1$, we are back to

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

- When $D = 2$ and writing

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

where $\rho = \text{corr}(X_1, X_2) \in [-1, 1]$ we have

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \times \exp\left(-\frac{1}{2(1-\rho^2)} \left\{ \frac{(x_1-\mu_1)^2}{\sigma_1^2} + \frac{(x_2-\mu_2)^2}{\sigma_2^2} - \frac{2(x_1-\mu_1)(x_2-\mu_2)}{\sigma_1\sigma_2} \right\}\right)$$

Graphical Illustrations

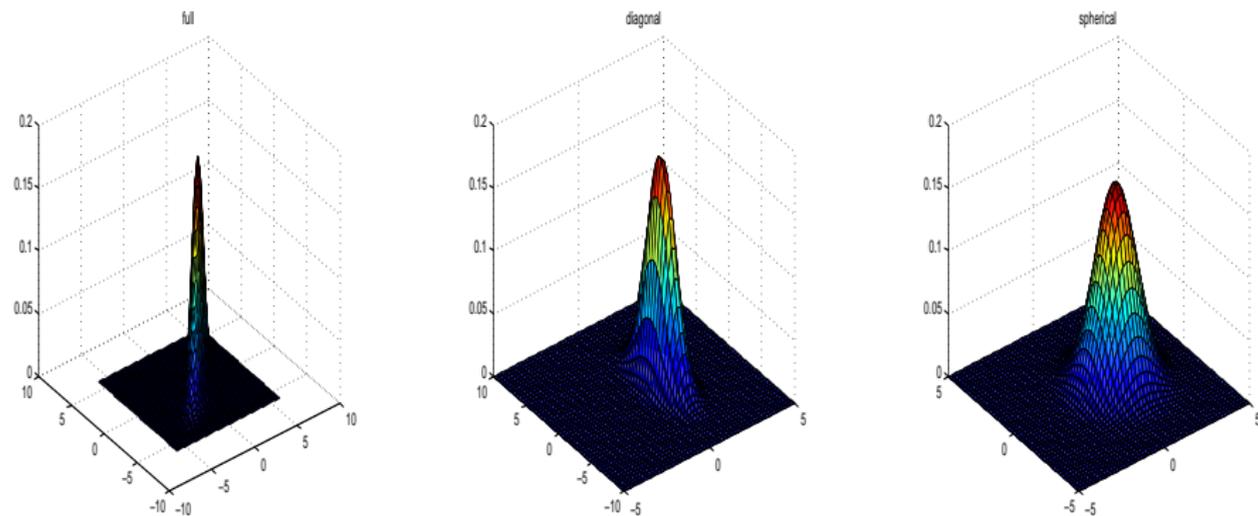


Illustration of 2D Gaussian pdfs for different covariance matrices (left): full, (middle): diagonal, (right): spherical.

Why are the contours of a multivariate Gaussian elliptical

- If we plot the values of \mathbf{x} s.t. $p(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$ is equal to a constant, i.e. s.t. $(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = c > 0$ where c is given, then we obtain an ellipse.
- Σ is a positive definite matrix so we have

$$\Sigma = U\Lambda U^\top$$

where U is an orthonormal matrix of eigenvectors, i.e. $U^\top U = I$, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_D)$ with $\lambda_k \geq 0$ is the diagonal matrix of eigenvalues.

- Hence, we have

$$\Sigma^{-1} = \left(U^\top\right)^{-1} \Lambda^{-1} U^{-1} = U\Lambda^{-1}U^\top = \sum_{k=1}^D \frac{1}{\lambda_k} \mathbf{u}_k \mathbf{u}_k^\top,$$

so

$$(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{k=1}^D \frac{y_k^2}{\lambda_k} \text{ where } y_k = \mathbf{u}_k^\top (\mathbf{x} - \boldsymbol{\mu}).$$

Graphical Illustrations

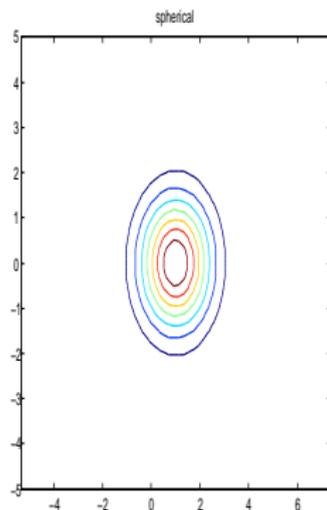
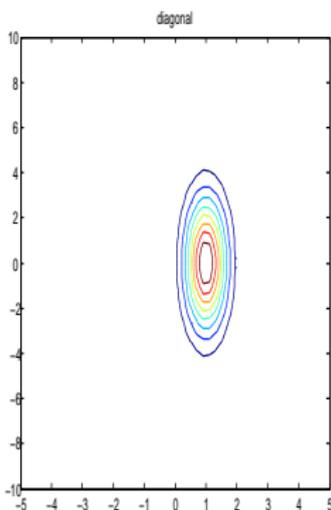
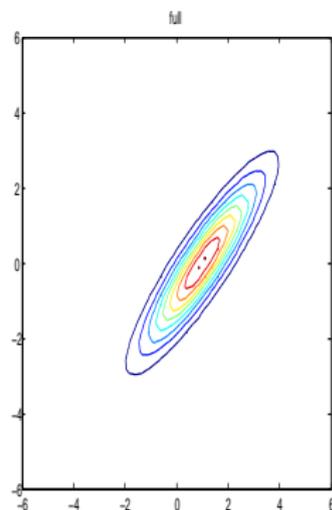


Illustration of 2D Gaussian pdfs level sets for different covariance matrices (left): full, (middle): diagonal, (right): spherical.

Properties of Multivariate Gaussians

- Marginalization is straightforward.
- Conditioning is easy; e.g. if $\mathbf{X} = (\mathbf{X}_1 \ \mathbf{X}_2)$ with

$$p(\mathbf{x}) = p(\mathbf{x}_1, \mathbf{x}_2) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$$

then

$$p(\mathbf{x}_1 | \mathbf{x}_2) = \frac{p(\mathbf{x}_1, \mathbf{x}_2)}{p(\mathbf{x}_2)} = \mathcal{N}(\mathbf{x}_1; \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

with

$$\begin{aligned} \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2), \\ \boldsymbol{\Sigma}_{1|2} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}. \end{aligned}$$

Independence and Correlation for Gaussian Variables

- It is well-known that independence implies uncorrelations; i.e. if the components (X_1, \dots, X_D) of a vector \mathbf{X} are independent then they are uncorrelated. However, uncorrelated does not imply independence in the general case.
- If the components (X_1, \dots, X_D) of a vector \mathbf{X} distributed according to a multivariate Gaussian are uncorrelated then they are independent.
- **Proof.** If (X_1, \dots, X_D) are uncorrelated then $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_D^2)$

and $|\Sigma| = \prod_{k=1}^D \sigma_k^2$ so

$$p(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \prod_{k=1}^D p(x_k | \mu_k, \sigma_k^2) = \prod_{k=1}^D \mathcal{N}(x_k; \mu_k, \sigma_k^2)$$

ML Parameter Learning for Multivariate Gaussian

- Consider data $\{\mathbf{x}^i\}_{i=1}^N$ where $\mathbf{x}^i \in \mathbb{R}^D$ and assume they are independent and identically distributed from $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$.
- The ML parameter estimates of $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ maximize by definition

$$\begin{aligned} & \sum_{i=1}^N \log \mathcal{N}(\mathbf{x}^i; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ = & -\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}^i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}^i - \boldsymbol{\mu}). \end{aligned}$$

- We obtain after painful calculations the fairly intuitive results

$$\hat{\boldsymbol{\mu}} = \frac{\sum_{i=1}^N \mathbf{x}^i}{N}, \quad \hat{\boldsymbol{\Sigma}} = \frac{\sum_{i=1}^N (\mathbf{x}^i - \hat{\boldsymbol{\mu}}) (\mathbf{x}^i - \hat{\boldsymbol{\mu}})^\top}{N}.$$

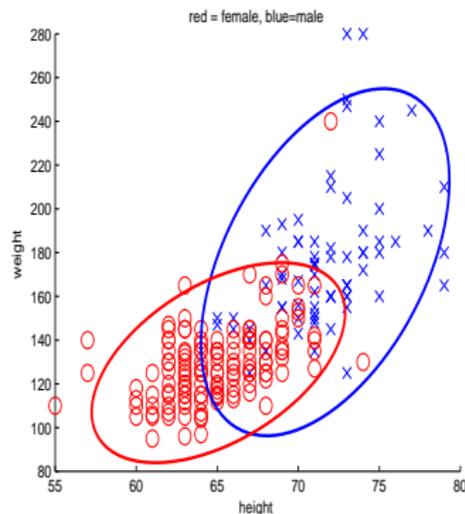
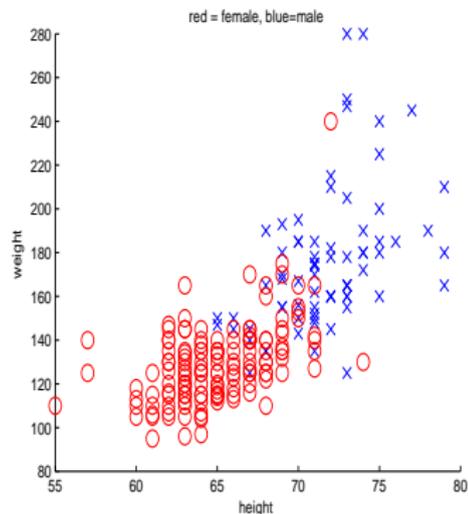
Application to Supervised Learning using Bayes Classifier

- Assume you are given some training data $\{\mathbf{x}^i, y^i\}_{i=1}^N$ where $\mathbf{x}^i \in \mathbb{R}^D$ and $y^i \in \{1, 2, \dots, C\}$ can take C different values.
- Given an input test data \mathbf{x} , you want to predict/estimate the output y associated to \mathbf{x} .
- Previously we have followed a probabilistic approach

$$p(y = c | \mathbf{x}) = \frac{p(y = c) p(\mathbf{x} | y = c)}{\sum_{c'=1}^C p(y = c') p(\mathbf{x} | y = c')}.$$

- This requires modelling and learning the parameters of the class conditional density of features $p(\mathbf{x} | y = c)$.

Height Weight Data



(left) Height/Weight data for female/male (right) 2d Gaussians fit to each class. 95% of the proba mass is inside the ellipse

Supervised Learning using Bayes Classifier

- Assume we pick

$$p(\mathbf{x}|y=c) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$$

and $p(y=c) = \pi_c$ then

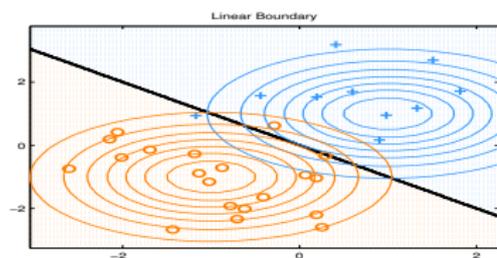
$$\begin{aligned} p(y=c|\mathbf{x}) &\propto \pi_c |\boldsymbol{\Sigma}_c|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}_c^{-1}(\mathbf{x} - \boldsymbol{\mu}_c)\right) \\ &= \exp\left(\boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}_c^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}_c^{-1} \boldsymbol{\mu}_c + \log \pi_c\right) \exp\left(-\frac{1}{2} \mathbf{x}^\top \boldsymbol{\Sigma}_c^{-1} \mathbf{x}\right) \end{aligned}$$

- For models where $\boldsymbol{\Sigma}_c = \boldsymbol{\Sigma}$ then this is known as *linear discriminant analysis*

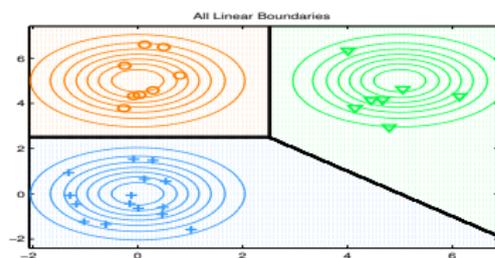
$$p(y=c|\mathbf{x}) = \frac{\exp\left(\beta_c^\top \mathbf{x} + \gamma_c\right)}{\sum_{c'=1}^C \exp\left(\beta_{c'}^\top \mathbf{x} + \gamma_{c'}\right)}$$

where $\beta_c = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c$, $\gamma_c = -\frac{1}{2} \boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c$ and the model is very similar to logistic regression.

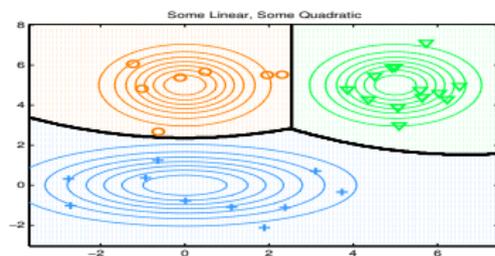
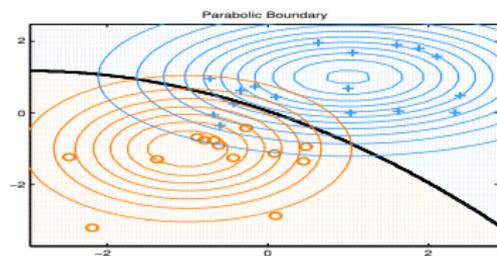
Decision Boundaries



(a)



(b)



Decision boundaries in 2D for 2 and 3 class case.

Binary Classification

- Consider the case where $C = 2$ then one can check that

$$p(y = 1 | \mathbf{x}) = g\left((\beta_1 - \beta_0)^\top \mathbf{x} + \gamma_1 - \gamma_0\right)$$

- We have

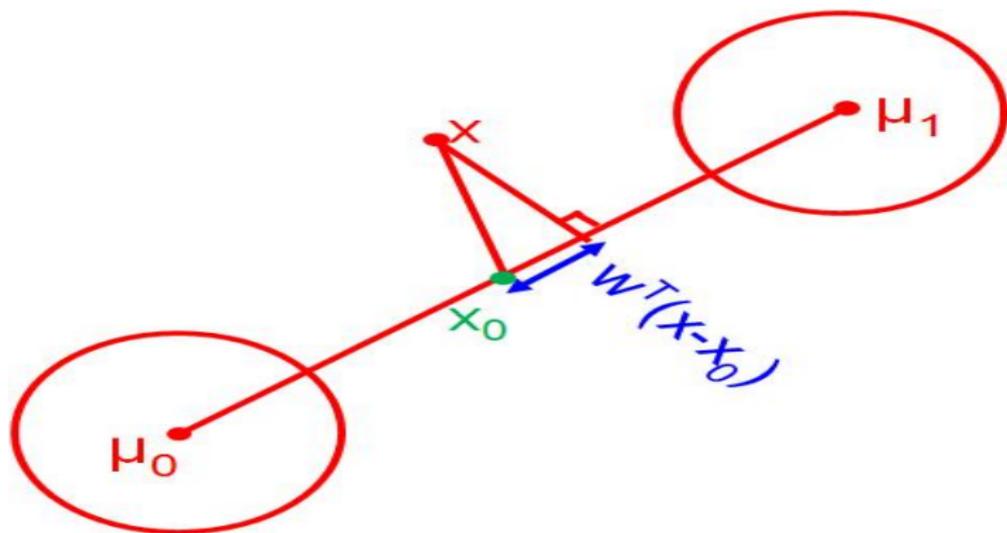
$$\begin{aligned}\gamma_1 - \gamma_0 &= -\frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) + \log(\pi_1/\pi_0), \\ \mathbf{x}_0 &: = \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) - \frac{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \log(\pi_1/\pi_0)}{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)} \\ \mathbf{w} &: = \beta_1 - \beta_0 = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)\end{aligned}$$

then

$$p(y = 1 | \mathbf{x}) = g\left(\mathbf{w}^\top (\mathbf{x} - \mathbf{x}_0)\right)$$

- \mathbf{x} is shifted by \mathbf{x}_0 and then projected onto the line \mathbf{w} .

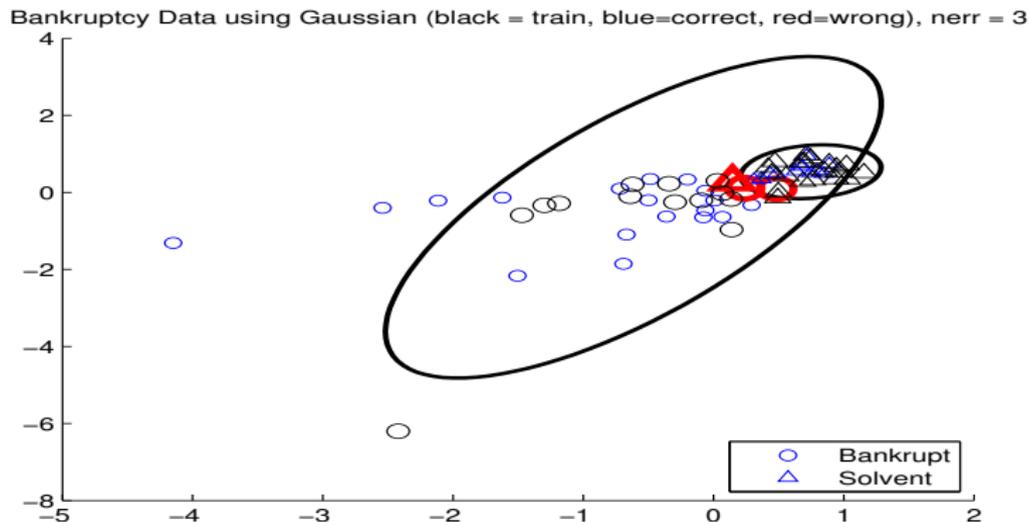
Binary Classification



Example where $\Sigma = \sigma^2 I$ so \mathbf{w} is in the direction of $\mu_1 - \mu_0$.

- When the model for class conditional densities is correct, resp. incorrect, generative classifiers will typically outperform, resp. underperform, discriminative classifiers for large enough datasets.
- Generative classifiers can be difficult to learn whereas Discriminative classifiers try to learn directly the posterior probability of interest.
- Generative classifiers can handle missing data easily, discriminative methods cannot.
- Discriminative can be more flexible; e.g. substitute \mathbf{x} to $\Phi(\mathbf{x})$.

Application



Discriminant analysis on the bankruptcy data set. Gaussian class conditional densities. Estimated labels, based on the posterior proba of belonging to each class, are computed. If incorrect, the point is colored read, otherwise in blue (Training data are black).