

$$p(\underline{x}, \underline{y}) = N\left(\begin{bmatrix} \underline{\mu}_x \\ \underline{\mu}_y \end{bmatrix}; \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{xy}^T & \Sigma_{yy} \end{bmatrix}\right)$$

$$\Rightarrow p(\underline{x}) = N(\underline{\mu}_x, \Sigma_{xx}); p(\underline{y}) = N(\underline{\mu}_y, \Sigma_{yy})$$

$$p(\underline{x} | \underline{y}) = N(\underline{\mu}_x + \Sigma_{xy} \Sigma_{yy}^{-1} (\underline{y} - \underline{\mu}_y), \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T)$$

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$

$$H(X) = E_{p(\underline{x})}[-\ln p(\underline{x})] = \frac{1}{2} [n_x \ln(2\pi e) + \ln |\Sigma_{xx}|]; H(Y) = \frac{1}{2} [n_y \ln(2\pi e) + \ln |\Sigma_{yy}|]$$

$$\begin{aligned} H(X, Y) &= \frac{1}{2} [(n_x + n_y) \ln(2\pi e) + \ln |\Sigma_{xx}| + \ln |\Sigma_{yy} - \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}|] \\ &= \frac{1}{2} [(n_x + n_y) \ln(2\pi e) + \ln |\Sigma_{yy}| + \ln |\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T|] \end{aligned}$$

$$I(X; Y) = \frac{1}{2} [\ln |\Sigma_{xx}| - \ln |\Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{xy}^T|] = \frac{1}{2} [\ln |\Sigma_{yy}| - \ln |\Sigma_{yy} - \Sigma_{xy}^T \Sigma_{xx}^{-1} \Sigma_{xy}|]$$

$$x, y \text{ scalars} : \Sigma_{xx} = \sigma_x^2; \Sigma_{xy} = \rho_{xy} \sigma_x \sigma_y; \Sigma_{yy} = \sigma_y^2$$

$$I(X; Y) = \frac{1}{2} [\ln \sigma_x^2 - \ln(\sigma_x^2(1 - \rho_{xy}^2))] = -\frac{1}{2} \ln(1 - \rho_{xy}^2) = \ln \frac{1}{\sqrt{1 - \rho_{xy}^2}}$$

Why sigmoid and softmax?

$$P(z=1 | \underline{x}) = \frac{p(\underline{x} | z=1)P(z=1)}{p(\underline{x})} = \frac{p(\underline{x} | z=1)P(z=1)}{p(\underline{x} | z=1)P(z=1) + p(\underline{x} | z=0)P(z=0)}$$

$$= \frac{1}{1 + \exp(-\ln \frac{p(\underline{x} | z=1)}{p(\underline{x} | z=0)} - \ln \frac{P(z=1)}{P(z=0)})} = \frac{1}{1 + \exp(-h(\underline{x}, \underline{w}))}$$

$$\Rightarrow P(z=1 | \underline{x}) = \frac{\exp(f(\underline{x}, w_1))}{\exp(f(\underline{x}, w_1)) + \exp(f(\underline{x}, w_0))} = \frac{1}{1 + \exp(-[f(\underline{x}, w_1) - f(\underline{x}, w_0)])}$$

$$P(z=i | \underline{x}) = \frac{\exp(f(\underline{x}, w_i))}{\sum_{j=1}^C \exp(f(\underline{x}, w_j))} = s_i(\underline{x}, \underline{w}) \dots \text{soft max}$$

$$\nabla_{w_i} s_i(\underline{x}, \underline{w}) = s_i(\underline{x}, \underline{w})(1 - s_i(\underline{x}, \underline{w})); \nabla_{w_j} s_i(\underline{x}, \underline{w}) = -s_i(\underline{x}, \underline{w})s_j(\underline{x}, \underline{w}); j \neq i$$

$$\begin{aligned} -\ln P(z_n | \underline{x}_n, \underline{w}) &= -z_n \underbrace{\ln P(z_n=1 | \underline{x}_n, \underline{w})}_{g(y_n)} - (1 - z_n) \underbrace{\ln P(z_n=0 | \underline{x}_n, \underline{w})}_{1-g(y_n)} \\ &= -z_n \ln g(y_n) - (1 - z_n) \ln(1 - g(y_n)) \end{aligned}$$

Data Distribution: $p = z_n \in \{0,1\}$; Classifier output: $q = \{g(y_n), 1 - g(y_n)\}$

$$\begin{aligned}
 \text{Cross-entropy } H_n(p, q) &= -z_n \ln g(y_n) - (1 - z_n) \ln(1 - g(y_n)) \\
 &= -z_n \ln z_n - (1 - z_n) \ln(1 - z_n) \\
 &\quad + [z_n \ln \frac{z_n}{g(y_n)} + (1 - z_n) \ln \frac{(1 - z_n)}{(1 - g(y_n))}] \\
 &= H_n(p) + KL(p \parallel q)
 \end{aligned}$$

$$L(x) = L(f(h(g(x))))$$

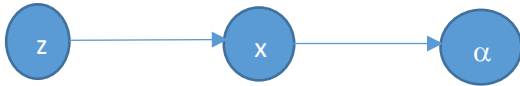
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial h} \frac{\partial h}{\partial g} \frac{\partial g}{\partial x} = \lambda_g \frac{\partial g}{\partial x}$$

$$\lambda_f = \frac{\partial L}{\partial f}; \lambda_h = \frac{\partial L}{\partial h} = \lambda_f \frac{\partial f}{\partial h};$$

$$\lambda_g = \frac{\partial L}{\partial g} = \lambda_h \frac{\partial h}{\partial g}$$

$$x \rightarrow g \rightarrow h \rightarrow f \rightarrow L$$

$$\frac{\partial L}{\partial x} = \lambda_g \frac{\partial g}{\partial x} \leftarrow \lambda_g = \frac{\partial L}{\partial g} = \lambda_h \frac{\partial h}{\partial g} \leftarrow \lambda_h = \frac{\partial L}{\partial h} = \lambda_f \frac{\partial f}{\partial h} \leftarrow \lambda_f = \frac{\partial L}{\partial f}$$



$$p(\alpha, \underline{x}, z) = P(z)P(\underline{x}/z)P(\alpha/\underline{x})$$

$$ECM = \sum_{j=1}^C \sum_{i=0}^C \lambda_{ij} P(\alpha = i, z = j)$$

$$= \int_{\underline{x}} \sum_{i=0}^C P(\alpha = i/\underline{x}) \cdot \left[\sum_{j=1}^C \lambda_{ij} p(\underline{x}/z = j) \cdot P(z = j) \right] d\underline{x}$$

$$\Rightarrow \text{Pick action } \alpha = k \text{ (class } \hat{z} = k \text{), if } k = \arg \min_{i \in \{0,1,2,\dots,C\}} \sum_{j=1}^C \lambda_{ij} p(\underline{x}/z = j) P(z = j)$$

$$= \arg \min_{i \in \{0,1,2,\dots,C\}} \sum_{j=1}^C \lambda_{ij} p(z = j/\underline{x})$$

$$\Lambda = \begin{bmatrix} \lambda_{01} & \lambda_{02} \\ \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{bmatrix}; \lambda_{11} < \lambda_{12}; \lambda_{22} < \lambda_{21} \Rightarrow \text{correct decisions have less cost}$$

Since $P(z = 2/\underline{x}) = 1 - P(z = 1/\underline{x})$, the decision rule is:

$$k = \arg \min_{i \in \{0,1,2\}} \{(\lambda_{i1} - \lambda_{i2})P(z = 1/\underline{x}) + \lambda_{i2}\}$$

Special case 1: $\lambda_{01} = \lambda_{02} = \lambda_r = \infty \Rightarrow$ Likelihood Ratio Rule

$$\frac{p(\underline{x}/z = 1)}{p(\underline{x}/z = 2)} \geq \frac{(\lambda_{12} - \lambda_{22})P(z = 2)}{(\lambda_{21} - \lambda_{11})P(z = 1)} \Rightarrow \hat{z} = 1; \text{otherwise } \hat{z} = 2 \text{ (Prove it)}$$

Special case 2: $\lambda_{01} = \lambda_{02} = \lambda_r$

Reject range for $P(z = 1/\underline{x})$ exists if $\frac{\lambda_r - \lambda_{12}}{\lambda_{11} - \lambda_{12}} > \frac{\lambda_r - \lambda_{22}}{\lambda_{21} - \lambda_{22}}$; else no reject decision (Prove it)

Special case 3: $\lambda_{01} = \lambda_{02} = \lambda_r; \lambda_{11} = \lambda_{22} = 0; \lambda_{12} > 0; \lambda_{21} > 0$

Reject range for $P(z = 1/\underline{x})$ exists if $\frac{1}{\lambda_r} > \frac{1}{\lambda_{12}} + \frac{1}{\lambda_{21}}$; else no reject decision (Prove it)

Pick action $\alpha = k, k \in \{0,1,2,\dots,C\}$

If $k = \arg \min \{\lambda_r, \lambda_e[1 - P(z = 1/\underline{x})], \dots, \lambda_e[1 - P(z = C/\underline{x})]\}$

$$= \arg \min \left\{ \frac{\lambda_r}{\lambda_e}, [1 - P(z = 1/\underline{x})], \dots, [1 - P(z = C/\underline{x})] \right\} \because \lambda_e > 0$$

$$= \arg \min \left\{ \frac{\lambda_r}{\lambda_e} - 1, -P(z = 1/\underline{x}), \dots, -P(z = C/\underline{x}) \right\}$$

$$= \arg \max \left\{ \underbrace{1 - \frac{\lambda_r}{\lambda_e}}_{\beta}, P(z = 1/\underline{x}), \dots, P(z = C/\underline{x}) \right\}$$

Decide for $\alpha = k$ (class $\hat{z} = k$): $\begin{cases} \text{if } P(z = k/\underline{x}) = \max_{j \in \{1,2,\dots,C\}} P(z = j/\underline{x}) \text{ and } P(z = k/\underline{x}) > \beta \\ \text{otherwise} & \text{reject (action } \alpha = 0) \end{cases}$

$$\text{MAP: } \beta = 0 \Rightarrow \lambda_r = \lambda_e \Rightarrow \arg \max_{j \in \{1,2,\dots,C\}} P(z = j/\underline{x})$$

ML: If $p(z = j) = \frac{1}{C} \forall j \Rightarrow \arg \max_{j \in \{1, 2, \dots, C\}} \ln p(\underline{x} / z = j)$

Properties of Gaussians: Sums, Linear Transformations, Conditionals

$$\underline{x} = \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix} \text{ and } \underline{x} \sim N(\underline{\mu}, \Sigma); \underline{\mu} = \begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{bmatrix}; \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix}; J = \Sigma^{-1} = \begin{bmatrix} J_{11} & J_{12} \\ J_{12}^T & J_{22} \end{bmatrix}$$

$$J_{11} = (\Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T)^{-1}; J_{22} = (\Sigma_{22} - \Sigma_{12}^T \Sigma_{11}^{-1} \Sigma_{12})^{-1}; J_{12} = -(\Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T)^{-1} \Sigma_{12} \Sigma_{22}^{-1}$$

Then, the marginal & conditional densities are also Gaussian

$$p(\underline{x}_2) = N(\underline{\mu}_2, \Sigma_{22})$$

$$p(\underline{x}_1 | \underline{x}_2) = N(\underline{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2), \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T) \\ = N(\underline{\mu}_1 - J_{11}^{-1} J_{12} (\underline{x}_2 - \underline{\mu}_2), J_{11}^{-1})$$

- Covariance matrix captures marginal independencies between variables

$$\Sigma_{ij} = 0 \Leftrightarrow x_i \text{ \& } x_j \text{ are independent (or) } x_i \perp x_j$$

- Information matrix $J = \Sigma^{-1}$ captures conditional independencies

$$J_{ij} = 0 \Leftrightarrow x_i \perp x_j \mid \{ \{x_1, x_2, \dots, x_n\} - \{x_i, x_j\} \}$$

Non-zero entries in J correspond to edges in the dependency network

- Simulating multivariate Gaussian random variables
- Discriminants

$$\max g_i(\underline{x}) = -\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(z = i)$$

$$\min \frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \frac{1}{2} \ln |\Sigma_i| - \ln P(z = i)$$

$$g_i(\underline{x}) = -\frac{1}{2} (\underline{x} - \underline{\mu}_i)^T \Sigma_i^{-1} (\underline{x} - \underline{\mu}_i) + \ln P(z = i) \\ = -\frac{1}{2} \underline{x}^T \Sigma_i^{-1} \underline{x} + \underline{\mu}_i^T \Sigma_i^{-1} \underline{x} - \left[\frac{1}{2} \underline{\mu}_i^T \Sigma_i^{-1} \underline{\mu}_i - \ln P(z = i) \right]$$

$$g_i(\underline{x}) = \underline{\mu}_i^T \Sigma_i^{-1} \underline{x} - \left[\frac{1}{2} \underline{\mu}_i^T \Sigma_i^{-1} \underline{\mu}_i - \ln P(z = i) \right]$$

$$= \underline{w}_i^T \underline{x} - w_{i0} \quad \dots \text{linear rule}$$

sigmoid, soft max for posteriors

$$g_i(\underline{x}) = \frac{1}{\sigma^2} \underline{\mu}_i^T \underline{x} - \left(\frac{1}{2\sigma^2} \underline{\mu}_i^T \underline{\mu}_i - \ln P(z = i) \right)$$

$$= \underline{w}_i^T \underline{x} - w_{i0}$$

$$g(\underline{x}) = \sigma^2 [g_1(\underline{x}) - g_2(\underline{x})]$$

- $$= (\underline{\mu}_1 - \underline{\mu}_2)^T \underline{x} - \left[\frac{1}{2} (\underline{\mu}_1^T \underline{\mu}_1 - \underline{\mu}_2^T \underline{\mu}_2) - \sigma^2 \ln \frac{P(z=1)}{P(z=2)} \right]$$

$$= \underline{w}^T \underline{x} - w_o = \underline{w}^T (\underline{x} - \underline{x}_o) = 0; \underline{x}_o = \frac{\underline{w}}{\underline{w}^T \underline{w}} w_o$$

Chernoff Bound:

$$\begin{aligned}
 P(\text{error}) &= \int_{\underline{x}} P(\text{error} / \underline{x}) p(\underline{x}) d\underline{x} \\
 &= \int_{\underline{x}} \min[P(z=1 / \underline{x}), P(z=2 / \underline{x})] p(\underline{x}) d\underline{x} \\
 &\leq [P(z=1)]^\beta [P(z=2)]^{1-\beta} \int_{\underline{x}} [p(\underline{x} / z=1)]^\beta [p(\underline{x} / z=2)]^{1-\beta} d\underline{x}
 \end{aligned}$$

If $a \geq b$ then $(a/b)^\beta \geq 1$

$$\Rightarrow (a/b)^\beta b \geq b$$

$$\Rightarrow a^\beta b^{1-\beta} \geq b$$

$$P(\text{error}) \leq [P(z=1)]^\beta [P(z=2)]^{1-\beta} e^{-k(\beta)} = e^{-k(\beta) + \beta \ln P(z=1) + (1-\beta) \ln [1 - P(z=1)]}$$

Gaussian: where $k(\beta) = \frac{\beta(1-\beta)}{2} \left\| \underline{\mu}_2 - \underline{\mu}_1 \right\|_{\beta \Sigma_1 + (1-\beta) \Sigma_2}^2 + \frac{1}{2} \ln \frac{[\beta \Sigma_1 + (1-\beta) \Sigma_2]}{\Sigma_1^\beta \Sigma_2^{1-\beta}}$

$$P(\text{error}) \leq \sqrt{P(z=1)P(z=2)} e^{-k(1/2)}$$

Bhattacharyya:

$$where \quad k(1/2) = \frac{1}{8} \left\| \underline{\mu}_2 - \underline{\mu}_1 \right\|_{\frac{\Sigma_1 + \Sigma_2}{2}}^2 + \frac{1}{2} \ln \frac{2}{\sqrt{|\Sigma_1| |\Sigma_2|}}$$

Discuss ROC

Binary case ... write as graph

Missing features

Noisy features

- Generative versus Discriminative?
- Can we exploit & generalize the forms of discriminant functions? Can we estimate discriminants directly?
- Different types of learning (unsupervised, supervised, semi-supervised, RL)
- How to handle missing data?
- How do we select features x for best classification/regression accuracy?
- Can we exploit dependency structure among features?
- What happens if classes and features change dynamically?
- How do we validate data driven models? How do we select the best model?