CS 479/679 Pattern Recognition Sample Midterm Exam

1. [20 pts] True/False Questions – To get credit, you **must** give brief reasons for your answers.

T F The Bayes rule is always an optimum classification rule in the sense that it minimizes the average probability error.

<u>True</u> ONLY under the assumption that $p(x/\omega i)$ and $P(\omega i)$ have been modelled/estimated correctly; in this case, the Bayes rule is optimum in the sense that it minimizes the average probability error.

T F Adding more features will in general improve classification accuracy.

<u>False</u> since adding more features will in general deteriorate classification accuracy due to the <u>"curse" of dimensionality</u> (i.e., the number of training data required increases exponentially with the number of features).

T F A causal relationship typically exists between correlated events.

<u>False</u> since correlated events are not necessarily causally related (i.e., correlations between two events can be caused by a third factor, called confounder, which affects both events).

T F A linear discriminant function could optimally separate the data between two classes when the features are uncorrelated.

<u>False</u> since the covariance matrix (diagonal in this case) for each class might be different (i.e., features can have different variances for each class). This is Case 3 (i.e., non-linear discriminant).

T F Given a set of N random variables $X=\{X_1, X_2, ..., X_N\}$, the pdf of any single random variable X_i can be computed from the joint pdf of X.

True by applying marginalization with respect to all other variables except X_i

- 2. [15 pts]] Short answer questions.
- a. Why do simpler models typically perform better than complex models in pattern recognition?

Complex models "memorize" the training data, therefore, they seneralize poorly.

b. What are the fundamental similarities/differences between Maximum A-Posteriori parameter estimation (MAP) and Bayesian parameter estimation (BE)?

P(0) Both MAP and BE discume known HAP: fixed point estimate BE: espinates a distribution

c. Under what conditions would the optimal decision boundary between two classes, each modeled by a Gaussian distribution, **not** pass from the midpoint of the line joining the means of the distributions?

When P(w,) + P(w) the decision boundary does not poss through the out apoint for Cours I, II (cose II could be equale or)

d. How could a matrix A be diagonalized? Are all matrices A diagonalizable? Explain.

PAP=1 P: columns are the eigenvectors of A A: diagonal, eigenalies of A - Not all matrices are diggenalizerable - Eigenvectors of A must be linearly independent

3. **[15 pts]** Consider a c-class classification problem; under what conditions would the optimal classifier be equivalent to (a) the minimum distance classifier? (b) the Mahalanobis distance classifier?

In both cases, we assume that each class is modelled by a multivariate Gaussian distribution.

(x) $\sum_{i=\sigma^2} T$ Assuaiing equal priors $\sum_{i=5}$ Kosyuning equal prims

4. **[15 pts]** Consider the following probability distribution:

$$p(x) = \begin{cases} (\theta+1)x^{\theta} & \text{for } 0 \le x \le 1\\ 0 & \text{otherwise} \end{cases}$$

Given n points x_1 , x_2 , ..., x_n sampled from the above distribution, derive a formula for the maximum likelihood estimate of ϑ .

$$P(D|\theta) = \prod_{i=1}^{n} p(X_i|\theta)$$

$$ln P(D|\theta) = \sum_{i=1}^{n} ln P(X_i|\theta) = \sum_{i=1}^{n} ln \left[(\theta_{i}) X_i^{\theta} \right] =$$

$$= \sum_{i=1}^{n} ln (\theta_{i}) + \sum_{i=1}^{n} ln (X_i^{\theta}) =$$

$$= ln (\theta_{i}) + \theta \sum_{i=1}^{n} ln (X_i^{\theta})$$

$$Take derivative and set cqual to zero:$$

$$\frac{\partial ln P(D|\theta)}{\partial \theta} = \frac{n}{\theta_{i}} + \sum_{i=1}^{n} ln (X_i) = 0$$

$$= ln (X_i) = \theta + i = -\frac{n}{\sum_{i=1}^{n} ln (X_i)}$$

$$= 0 = -1 - \frac{n}{\sum_{i=1}^{n} ln (X_i)}$$

5. **[20 pts]** Discuss the main ideas of **the parameter estimation** techniques presented in class and compare them with each other. Are their solutions always different? Explain.

Maximum Likelihood Estimation

It assumes that the parameters are fixed but unknown. It maximizes the likelihood $p(D/\theta)$ where D is the training data and θ the parameters to be estimated.

Maximum A-Posteriori Estimation:

It assumes that the parameters are fixed and have a known distribution $p(\theta)$. It maximizes the likelihood $p(D/\theta)p(\theta)$ where D is the training data and θ the parameters to be estimated.

Bayesian Estimation:

It assumes that the parameters are random variables with a known distribution $p(\theta)$. Using the data D, it converts $p(\theta)$ to $p(\theta/D)$. It estimates a distribution:

$$p(\mathbf{x}/D) = \int p(\mathbf{x}/\boldsymbol{\theta}) p(\boldsymbol{\theta}/D) d\boldsymbol{\theta}$$

Comparison

MLE/MAP are simpler than BE and their solution is simpler to interpret. They also have lower computational requirements compared to BE. The BE solution might not be of the form (model) assumed.

The BE solution is different from the MLE solution when the dataset is not very large. In theory, the BE solution converges to the MLE solution as $n \rightarrow \infty$

6. **[15 pts]** Using the Bayesian network below, answer the following question: suppose all we know is that a fish is thin, has medium lightness, and was caught in south Atlantic; what season is it now, most likely?



We need to compute
$$P(a_i/d_2, C_2, b_2)$$

 $P(a_i/d_2, C_2, b_2) = \frac{P(a_i, d_2, c_2, b_2)}{P(d_2, C_2, b_2)} \sum_{a_j} P(a_i, b_2, X_j, c_2, d_2)$
 $P(a_i/d_2, C_2, b_2) = \sum_{j} P(a_j, b_2, X_j, (c_2, d_2)) = \sum_{j} P(a_j, b_2, X_j, (c_2, d_2)) + P(a_j, b_2, X_2, (c_2, d_2)) = \sum_{j} P(a_j)P(b_2)P(X_1|a_j, b_2)P(c_2, X_j)P(d_2/X_1) / P(a_j)P(b_2)P(X_2|a_j, b_2) + P(a_j, b_2)P(X_2|a_j, b_2) = \sum_{j} P(a_j)P(b_2)P(X_2|a_j, b_2) = 0.0136 + P(a_j, a_j, b_2)P(a_j, b_2)P(a_j, b_2)P(a_j, b_2)P(a_j, b_2)P(a_j, b_2) = 0.0136 + P(a_j, a_j, b_2) = 0.0136 + P(a_j, a_j, b_2) = 0.0126 + P(a_j, a_j, b_2) = 0.0126 + P(a_j, b_2)P(a_j, b_2)P(a_j, b_2) = 0.0126 + P(a_j, b_2)P(a_j, b_2)P(a_j, b_2) = 0.0126 + P(a_j, b_2)P(a_j, b_2)P(a_j, b_2)P(a_j, b_2) = 0.0126 + P(a_j, b_2)P(a_j, b$