

Attempt all questions. Provide answers to all questions in a single page, neatly numbered in order. You may attach two further pages with any workings where useful, clearly numbered and in order of the questions. These additional pages will be looked at if the question requires any derivation and the answer you provided is incorrect.

Note the questions will have at least one correct answer. Where the question has more than one correct answer, you must select all the correct ones. For these questions, partial credit will usually not be given.

You should upload a maximum of three pages as a single pdf file.

Each question is worth *five marks*.

$P(C_2|X)$

Date TBD

$P(C_1|X)$

Timing Four hours within a 24 hour period

高斯属性判别分析 RRML: 15 次第

- To solve a two-class classification problem with classes denoted  $\mathcal{A}$  and  $\mathcal{B}$ , a classifier is designed to compute  $y = \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$  on input features  $\mathbf{x} \in \mathcal{R}^p$ , with  $\mathbf{A} \in \mathcal{R}^{p \times p}$ ,  $\mathbf{b} \in \mathcal{R}^p$  and  $c \in \mathcal{R}$  as its parameters. Once the parameters are determined, a decision rule of the form  $y > \theta \Rightarrow \mathcal{A}$  is applied. Which of the following conditions are necessary for the above classifier to be optimal.

$$P(C_1|X) \geq P(C_2|X) \Rightarrow \text{决策加权}$$

- The prior probabilities of the two classes are different: i.e.  $p(\mathcal{A}) \neq p(\mathcal{B})$ . 无法确定
- The class conditional densities are *mixture Gaussian distributions*: i.e.  $p(\mathbf{x}|\mathcal{A}) = \sum_{k=1}^{K_A} \lambda_k \mathcal{N}(\mathbf{m}_k, \Sigma_k)$  and likewise for class  $\mathcal{B}$ . 只影响偏移量

$$P(C_1|X) = \frac{1}{1+e^{-x}} \Rightarrow \text{决策加权}$$

- The class conditional probabilities are both Gaussian: i.e.  $p(\mathbf{x}|\mathcal{A}) = \mathcal{N}(\mathbf{m}_{\mathcal{A}}, \Sigma_{\mathcal{A}})$  and  $p(\mathbf{x}|\mathcal{B}) = \mathcal{N}(\mathbf{m}_{\mathcal{B}}, \Sigma_{\mathcal{B}})$ , where  $\Sigma_{\mathcal{A}} \neq \Sigma_{\mathcal{B}}$ . A 和 B 不同

$$P(C_1|X) = \frac{1}{1+e^{-x}}$$

- The class conditional probabilities are both Gaussian and should have identical an isotropic covariance metrics: i.e.  $p(\mathbf{x}|\mathcal{A}) = \mathcal{N}(\mathbf{m}_{\mathcal{A}}, \Sigma_{\mathcal{A}})$  and  $p(\mathbf{x}|\mathcal{B}) = \mathcal{N}(\mathbf{m}_{\mathcal{B}}, \Sigma_{\mathcal{B}})$ , where  $\Sigma_{\mathcal{A}} = \Sigma_{\mathcal{B}} = \sigma^2 I$ . 同性协方差矩阵

$$\text{取決于} \alpha = \ln \frac{P(C_2)}{P(C_1)} + \ln \frac{P(X|C_2)}{P(X|C_1)}$$

- A variable of interest,  $x$ , has an exponential distribution:

$$\log p(x|\theta) = \begin{cases} \theta \exp(-\theta x) & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Given  $N$  samples,  $x_1, x_2, \dots, x_N$ , drawn independently from  $P(x|\theta)$ , derive the maximum likelihood estimate of  $\theta$ .  $\hat{\theta} = \arg \max_{\theta} \log P(x|\theta)$  MLE  $P(x|\theta) = \prod_{i=1}^N P(x_i|\theta) = \prod_{i=1}^N \theta \exp(-\theta x_i)$

Your answer is:

$$\checkmark \quad \hat{\theta} = \frac{1}{\frac{1}{N} \sum_{n=1}^N x_n}$$

$$= \arg \max \sum_{i=1}^N \log P(x_i|\theta)$$

$$= -\sum_{i=1}^N \theta x_i + N \ln \theta$$

$$2. \quad \hat{\theta} = \prod_{n=1}^N x_n$$

$$= \arg \max \sum_{i=1}^N \log \theta \exp(-\theta x_i)$$

$$\text{取 } \ln P(x|\theta) = \sum_{i=1}^N \ln \theta \exp(-\theta x_i)$$

$$3. \quad \hat{\theta} = \frac{N!}{2\pi} \sum_{n=1}^N N x_n$$

$$= \arg \max \sum_{i=1}^N \log \theta - \frac{1}{2\pi} \theta x_i$$

$$= \sum_{i=1}^N (\ln \theta + (-\theta x_i)) = -\sum_{i=1}^N \theta x_i + N \ln \theta$$

$$4. \quad \hat{\theta} = \frac{1}{\sqrt{2\pi}} \sum_{n=1}^N x_n$$

$$= \arg \max$$

$$\text{取最大，对 } \ln P(x|\theta) \text{ 做导数}$$

$$\Rightarrow \frac{N}{\theta} = \frac{N}{\sqrt{2\pi}} \sum_{i=1}^N x_i \Rightarrow \theta = \frac{N}{\sum_{i=1}^N x_i} = \frac{1}{\frac{1}{N} \sum_{i=1}^N x_i}$$

$$\frac{\partial \ln P(x|\theta)}{\partial \theta} = -\frac{N}{N} + \frac{N}{\theta} = 0$$

## 解析题：

3. You are tasked with predicting the market price of an asset using its past values and several variables relating to the underlying economy within which the business operates. The dataset given to you spans six months of daily trading (125 items) of 100 variables. You are required to split the data into a training set and an evaluation set of equal sizes and use a linear model as predictor. You attempt to solve the problem of estimating regression coefficients by

回归系数

$$\underline{w = (X^T X + \gamma I)^{-1} X^T t} \rightarrow \text{最小二乘估计}$$

Which of the following is/are true?

2.  $w = (X^T X)^{-1} X^T t \dots \text{①}$

$x = [N \times M]$  的矩阵

$N < M$ , 过拟合

此时  $(X^T X)^{-1}$  不可逆

② 防止： $L_1 : P(w) = \|w\|_1$

$L_2 : P(w) = \|w\|_2 = \sqrt{w^T w}$  正则化

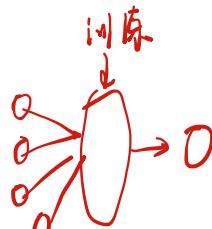
4. A researcher attempting to predict median house prices in different parts of Southampton from 15 covariates recommended by a group of estate agents, decides to use a multi-layer perceptron as defined below in the sklearn package. The researcher has access to a total of 200 data points and chooses a random subset of 150 for training, reserving the remaining 50 for a final evaluation of the model.

```
MLPClassifier(hidden_layer_sizes=1000,
               activation='relu', *,
               solver='adam',
               alpha=0.0001,
               batch_size='auto',
               learning_rate='constant',
               learning_rate_init=0.001,
               power_t=0.5,
               max_iter=200000 → 调小
               shuffle=True,
               random_state=None,
               tol=0.0001,
               verbose=False,
               warm_start=False,
               momentum=0.9,
               nesterovs_momentum=True,
               early_stopping=False → 设为 True
               validation_fraction=0.1,
               beta_1=0.9,
               beta_2=0.999,
               epsilon=1e-08,
               n_iter_no_change=10,
               max_fun=15000)
```

data ↴

当误差到一定时  
就停止，可以  
防止过拟合

150个训练 ↴ 极  
50个 test ↴ 有可能且拟合  
hidden 层有 1000 个 ↴ 太大了  
15 个维合



神经网络的图解



Which of the following statements is true?

1. It is possible that the error on the training set reduces to a negligibly small value.
2. The average error on the 50 test samples is likely to be disappointingly high.

过拟合

# 避免 over-fitting 的方法之一

3. If the model does not generalize, additional hidden layers should be used.
4. I would advice against setting early\_stopping=False 可以平衡过拟合
5. Your supervisor has asked you to apply Expectation-Maximisation for clustering. Which of the following statements is true?   
 1. The objective is guaranteed to converge to a global optimum.   
 2. Convergence rate depends on the properties of the likelihood function.   
 3. The log-likelihood,  $\ln p_{\theta}(\mathbf{X})$  decreases monotonically.   
 6. Consider Expectation-Maximisation (EM) for Gaussian Mixture Models (GMMs). Assume the GMM has  $K$  components and that all GMM components share the same covariance  $\Sigma_1 = \dots = \Sigma_K = \epsilon \mathbf{I}$ , where  $\mathbf{I}$  denotes the identity matrix and  $\epsilon$  is a constant. Which of the following statements is true?
1. EM approaches the  $k$ -means solution if  $\epsilon = 1$ .   
 2. EM approaches the  $k$ -means solution if  $\epsilon \rightarrow \infty$ .   
 3. EM approaches the  $k$ -means solution if  $\epsilon \rightarrow 0$ .

PRML

GMMs 是将每个样本  $\mathbf{x}$  按概率 软分配到  $K$  个高斯分布上。

$k$ -means 是将样本  $\mathbf{x}$  硬分布到  $K$  个点上。

$K$  个点可以看成  $\epsilon \rightarrow 0$  的高斯分布。



$k$ -Means : 硬分类  
 100% 放在某个类里  
 个体 (距离中心点的远近)  $\rightarrow$  只是干点  
 聚类 / 分类 (已知)

GMM : 软分类  
 $p(c_1)$  与  $p(c_2)$  的大小区别  
 存在一定方差

