# Mock Multiple Choice Questions — com4509/com6509

1. [probability] A random variable $X$ takes a value $x$, drawn from a set $X$ containing possible outcomes (or events). Let $p(x)$ denote the probability density function of $X$. The discrete type expected value of a random variable $X$ is defined as

a) $E[X] = \sum_{x \in X} \frac{d}{dx} p(x)$

b) $E[X] = \sum_{x \in X} p(x)$

**c)** $E[X] = \sum_{x \in X} x p(x)$

d) $E[X] = \sum_{x \in X} x^2 p(x)$

e) non of the above

**Expected Value**

The **expected value** (or mean, average) of a random variable $X$ is

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} x p(x) dx$$

- discrete type is $\mathbb{E}[X] = \sum_{x \in X} x p(x)$ for all possible events $X$

The expected value of a function $f(x)$ is

$$\mathbb{E}[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx$$

2. [probability] A random variable $X$ (or $Y$) takes a value $x$ (or $y$), drawn from a set $X$ (or $Y$) containing possible outcomes. Which one of the following is the correctly formulated?

Slide 1

a) $p(x|y) = p(y|x)p(x)$

**b)** $p(x|y) = \dfrac{p(y|x)p(x)}{\sum_{x \in X} p(y|x)p(x)}$

c) $p(x|y) = p(y|x)p(y)$

d) $p(x|y) = \dfrac{p(y|x)p(y)}{\sum_{x \in X} p(y|x)p(y)}$

e) $p(x|y) = \dfrac{p(y|x)}{\sum_{x \in X} p(y|x)}$

**Product Rule and Sum Rule**

The product rule of probability:

$$\underbrace{p(x, y)}_{\text{joint probability}} = \underbrace{p(y|x)}_{\text{conditional probability}} \cdot p(x)$$

The sum rule of probability:

$$\underbrace{p(y)}_{\text{marginal probability}} = \sum_{x \in X} p(x, y) = \sum_{x \in X} p(y|x)p(x)$$

**Bayes' Theorem**

Bayes' theorem immediately follows the product rule:

$$p(x|y) = \frac{p(x, y)}{p(y)} = \frac{p(y|x)p(x)}{\sum_{x \in X} p(y|x)p(x)}$$

3. [vector calculation] There are three vectors, $\mathbf{x} = (1, 3, -5)$, $\mathbf{y} = (4, -2, -1)$ and $\mathbf{z} = (-2, -5, 2)$. Which of the following is <u>correct</u>?

a) $\mathbf{x}$ and $\mathbf{y}$ are orthogonal.

**b)** $\mathbf{y}$ and $\mathbf{z}$ are orthogonal.   $\mathbf{y} \cdot \mathbf{z} = 0$

c) $\mathbf{x}$ and $\mathbf{z}$ are orthogonal.

d) $\mathbf{x}$ and $\mathbf{y}$ and $\mathbf{z}$ are orthogonal.

e) non of the above.

4. [multivariate normal distribution] Using an $n \times 1$ column vector $\mu$ and an $n \times n$ square matrix $\Sigma$, the multivariate normal distribution of an $n$-dimensional random vector $\mathbf{x}$ is

Slide 1

$$p(\mathbf{x}|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n \Sigma}} \exp\left\{ -\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu) \right\}$$

where $(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)$ is   ← $1 \times n$, $n \times n$, $n \times 1 = 1 \times 1$ (scalar value)

a) $n \times n$ square matrix

b) $n \times 1$ column vector

c) $1 \times n$ row vector

**d)** scalar value

e) non of the above

**Multivariate Gaussian (Normal) Distribution**

The probability density of the $k$-dimensional **Gaussian distribution** is

$$p(\mathbf{x} | \mu, \Sigma) = \frac{1}{\sqrt{2\pi^k |\Sigma|}} \exp\left( -\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu) \right)$$

where $\mu$ is the $k \times 1$ **mean vector** and $\Sigma$ is the $k \times k$ **covariance matrix**.

- $|\Sigma|$ and $\Sigma^{-1}$ are the **determinant** and the **inverse** of the covariance
- a symbol $^\top$ indicates the **transpose**

5. [overdetermined systems] Consider the following four overdetermined systems, each consisting of three equations and two unknowns $x$ and $y$:

$$\begin{cases} 2x+y=-1 \\ -3x+y=-2 \\ -x+y=1 \end{cases} \quad \begin{cases} 2x+y=1 \\ -2x+y=1 \\ -x+y=1 \end{cases} \quad \begin{cases} x-y=0 \\ 2x-2y=-2 \\ -x+y=1 \end{cases} \quad \begin{cases} 2x+y=1 \\ 2x-2y=-2 \\ -x+y=1 \end{cases}$$

$\times$ $\checkmark$ $\times$ $\checkmark$

How many systems above have solutions?

a) 0

b) 1

c) 2

d) 3

e) 4

6. [objective function] How many of the following four objective functions are linear in terms of their parameters?

$$f_1(x) = w_1 + w_2 x + w_3 x^2 \quad \checkmark$$
$$f_2(x) = w_1 + w_2 e^x + w_3 e^{x^2} \quad \checkmark$$
$$f_3(x) = w_1 + e^{w_2 x} + e^{w_3 x^2}$$
$$f_4(x) = w_1 + \log(x^{w_2}) + \log(x^{2w_2}) \quad \checkmark$$

a) 0

b) 1

c) 2

d) 3     Lab 4

e) 4

7. [gradient descent] Which one of the following is <u>incorrectly</u> described?

Lab 2
Slide 2

a) A stochastic gradient descent (SGD) algorithm must update parameters after calculating the gradient of the objective functions for all data points in a random order.

b) A stochastic gradient descent (SGD) algorithm can update parameters after calculating the gradient of the objective functions for a subset of all data points.

c) A stochastic gradient descent (SGD) algorithm can update parameters after calculating the gradient of the objective function for a single data point.

d) A standard gradient descent (GD) method can update parameters after calculating the gradient of the objective functions for all data points.

e) A standard gradient descent (GD) method can update parameters after calculating the gradient of the objective functions for all data points in a random order.

**Stochastic Gradient Descent Algorithm**

If the number of data points $n$ is small, gradient descent is fine, but sometimes (eg. 'Big Data') $n$ could be a billion.

• stochastic gradient descent is more similar to perceptron.
• it looks at gradient of one data point at a time rather than summing across *all* data points.
• this gives a stochastic estimate of gradient.

8. **[likelihood]** The likelihood of a parameter set $\theta$ for observed data $x$ is

**Likelihood**

a) $p(x|\theta)p(\theta)$

The likelihood of parameter values $\theta$ given data **x** is the probability for those observed data given those parameter values:

$$\mathcal{L}(\theta|x) = p(x|\theta)$$

Slide 3  b) $p(x|\theta)$

Slide 6  c) $\dfrac{p(x|\theta)p(\theta)}{p(x)}$

**Posterior Distribution**

Posterior distribution is found by combining the prior with the likelihood.

d) $p(\theta|x)$

- posterior distribution to represent our belief *after* we see the data of the likely value of the parameters
- the posterior being derived through **Bayes' rule**:

e) $p(\theta|x)p(x)$

$$p(c|y) = \frac{p(y|c)p(c)}{p(y)}$$

9. **[likelihood]** Which one of the following descriptions about a log likelihood is <u>correct</u>?

Slide 3  a) The logarithm is a monotonically ~~decreasing~~ *increasing* function, hence the log of a likelihood achieves its maximum at the point where the likelihood achieves its ~~minimum~~ *maximum*.

Slide 6  b) A posterior pdf (probability density function) is calculated as a **product** of the **log** likelihood and the prior pdf. $p(c|\mathbf{x}, \mathbf{y}, m, \sigma^2) = \frac{p(\mathbf{y}|\mathbf{x}, m, c, \sigma^2)p(c)}{p(\mathbf{y}|\mathbf{x}, m, \sigma^2)}$ (PDF)

c) A number of well known distributions, such as the normal distribution, have likelihood functions that contain products of factors involving exponentiation. The logarithm of such a function repllaces products with sums, hence easier to differentiate than the original function.

Slide 3  d) Probabilistic interpretation for an error function is a log likelihood of normal distribution. Thus maximising error function is equivalent to maximum likelihood with respect to parameters.

Slide 3  e) The log likelihood of the normal distribution involves exponentiation.

10. **[generalisation]** Which of the following describes cross validation <u>incorrectly</u>?

Lab 5  a) It assesses how the results of a statistical analysis will generalise to an unseen data set.

Slide 5  b) It tests the model using the validation dataset in order to reduce the chance of an overfitting problem.

c) It partitions data into complementary subsets for training and for validation, which is repeated using different partitions and finally the validation results are averaged.

d) Leave-one-out validation is its special case, where partition is made with one validation data point and the rest for training data.

e) It produces meaningful results if the validation set and training set are drawn from the totally different population.

11. **[unsupervised learning]** Which one of the following is <u>not</u> considered as an approach to unsupervised learning?

Slide 7

**Supervised Learning**

Supervised learning is learning where each data has a label.

- training example: set of pairs consisting of an input object and a desired output value
- recall the bias-variance tradeoff (**week 5**)

(eg) regression analysis

a) support vector machine (SVM)

b) k-means clustering

**Unsupervised Learning**

In unsupervised learning we have no labels for the data. It is often thought of as structure discovery, such as finding features in the data.

c) latent variable models

Unsupervised learning techniques:

d) factor analysis

- clustering (eg: k-means)
- latent variable models (eg: EM algorithm)
- blind signal separation (eg: PCA, ICA, SVD)

e) principal component analysis (PCA)

Questions 12 to 15 are about python code. The `numpy` library has been imported as `np` and we are given `u` and `v` as one dimensional numpy arrays. Note that `numpy.outer` calculates the outer product of two vectors.

12. [python code] Which one of the following pieces of python code calculates $\mathbf{Y} = \mathbf{uv}^\top$? Note that $\mathbf{u}$ and $\mathbf{v}$ are both column vectors.

    a) `y = u*v`

    b) `y = np.sum(u*v)`

    c) `y = np.sum(v*u**2)`

    d) `y = np.sum(u*v)**2`

    e) `y = np.outer(u,v)`

13. [python code] Which one of the following pieces of python code calculates $y = \mathbf{v}^\top \mathbf{uu}^\top \mathbf{v}$? Note that $\mathbf{u}$ and $\mathbf{v}$ are both column vectors.

    a) `y = u*v`

    b) `y = np.sum(u*v)`

    c) `y = np.sum(v*u**2)`

    d) `y = np.sum(u*v)**2`

    e) `y = np.outer(u,v)`

14. [python code] Which one of the following pieces of python code calculates $y = tr(\mathbf{uv}^\top)$? Note that $\mathbf{u}$ and $\mathbf{v}$ are both column vectors, and that $tr(\mathbf{A})$ indicates the sum of the diagonal elements of the square matrix $\mathbf{A}$.

    a) `y = u*v`

    b) `y = np.sum(u*v)`

    c) `y = np.sum(v*u**2)`

    d) `y = np.sum(u*v)**2`

    e) `y = np.outer(u,v)`

15. [python code] Which one of the following pieces of python code calculates $y = \mathbf{v}^\top diag(\mathbf{u})$? Note that $\mathbf{u}$ and $\mathbf{v}$ are both column vectors, and that $diag(\mathbf{z})$ forms a diagonal matrix with diagonal elements given by elements of $\mathbf{z}$.

    a) `y = u*v`   位乘

    b) `y = np.sum(u*v)`

    c) `y = np.sum(v*u**2)`

    d) `y = np.sum(u*v)**2`

    e) `y = np.outer(u,v)`