

## APPHYS 237 / BIO 251 – Mathematical Prerequisites

This course is designed to be accessible to students with a broad range of backgrounds in physics, biology, applied mathematics, engineering, etc. The main prerequisites for the course are mathematical, in particular, proficiency in multivariable calculus and basic familiarity with differential equations and probability theory.

The goal of this document is to list some of the mathematical concepts you should be familiar with in order to thrive in the course. Depending on your learning style, it may be possible to pick up one or two of these concepts during the course itself (particularly if you are the kind of person who likes learning through applications). However, we will be building on many of these concepts to introduce new mathematical tools that we will need for the course, so it may be difficult to catch up if you are unfamiliar with too many of the prerequisites.

### Mathematical Prerequisites

#### Calculus

- **Derivatives**, e.g. if  $f(x) = \exp(-x^2)$ , then

$$f'(x) \equiv \frac{\partial f}{\partial x} = -2x \exp(-x^2) \quad (1)$$

$$f''(x) \equiv \frac{\partial^2 f}{\partial x^2} = -2 \exp(-x^2) + 4x^2 \exp(-x^2) \quad (2)$$

- **Integrals**, e.g.

$$\int x e^{-x^2} dx = -\frac{1}{2} e^{-x^2} \quad (3)$$

$$\int x e^{-x} dx = -(1+x)e^{-x} \quad (4)$$

- **Series expansions**, e.g. Taylor series around  $x = a$ :

$$f(x) = \sum_{k=0}^{\infty} \frac{(x-a)^k}{k!} \left. \frac{\partial^k f}{\partial x^k} \right|_{x=a} \quad (5)$$

e.g., around  $x = 0$ ,

$$e^x = 1 + x + \frac{x^2}{2} + \dots = \sum_{k=0}^{\infty} \frac{x^k}{k!} \quad (6)$$

$$\frac{1}{1+x} = 1 - x + x^2 + \dots = \sum_{k=0}^{\infty} (-1)^k x^k \quad (7)$$

$$\log(1+x) = x - \frac{x^2}{2} + \dots = \sum_{k=1}^{\infty} (-1)^k \frac{x^k}{k} \quad (8)$$

Eq. 7 also yields a formula for the sum of a geometric series with a *finite* number of terms:

$$\sum_{k=0}^n x^k = \frac{1-x^{n+1}}{1-x} \quad (9)$$

Using truncated Taylor series as approximations, e.g. when “ $x$  is close to 0”,

$$e^{-x} \approx 1 - x \quad (10)$$

$$(1 + x)^k \approx 1 + kx \quad (11)$$

$$\log(1 + x) \approx x \quad (12)$$

## Differential equations

- Ordinary differential equations, e.g.

**Exponential growth:**  $\frac{\partial f}{\partial t} = sf \rightarrow f(t) = f(0)e^{st}$  (13)

**Logistic growth:**  $\frac{\partial f}{\partial t} = sf(1 - f) \rightarrow f(t) = \frac{f(0)e^{st}}{1 + f(0)(e^{st} - 1)}$  (14)

Verification of solutions by substitution. Finding solutions using **separation of variables** and/or **integrating factors**, e.g.

$$\frac{\partial f}{\partial t} = sf + U(t) \rightarrow f(0)e^{st} + e^{st} \int_0^t U(t')e^{-st'} dt' \quad (15)$$

- Partial differential equations, e.g.

**Advection equation:**  $\frac{\partial f}{\partial t} = -v \frac{\partial f}{\partial x} \rightarrow f(x, t) = f(x - vt, 0)$  (16)

**Diffusion equation:**  $\frac{\partial f}{\partial t} = D \frac{\partial^2 f}{\partial x^2} \rightarrow f(x, t) = \int_{-\infty}^{\infty} \frac{f(x', 0)}{\sqrt{4\pi Dt}} e^{-\frac{(x-x')^2}{4Dt}} dx'$  (17)

Verification of solutions by substitution. (Advanced) Obtaining solutions via **separation of variables** or **Fourier/Laplace transforms**

## Probability and statistics

- Discrete random variables,  $\Pr[a \leq X \leq b] = \sum_{k=a}^b \Pr[X = k]$ , e.g.

**Bernoulli( $p$ ) distribution:**  $\Pr[X = k] = \begin{cases} p & \text{if } k = 1 \\ 1 - p & \text{if } k = 0 \end{cases}, \quad k = 0, 1$  (18)

**Binomial( $n, p$ ) distribution:**  $\Pr[X = k] = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, \dots, n$  (19)

**Poisson( $\lambda$ ) distribution:**  $\Pr[X = k] = \frac{\lambda^k}{k!} e^{-\lambda}, \quad k = 0, 1, \dots, \infty$  (20)

- **Continuous random variables**,  $\Pr[a \leq X \leq b] = \int_a^b p(x) dx$ , e.g.

**Exponential( $\lambda$ ) distribution:**  $p(x)dx = \lambda e^{-\lambda x} dx$ ,  $x \geq 0$  (21)

**Gaussian( $\mu, \sigma^2$ ) distribution:**  $p(x)dx = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$ ,  $-\infty < x < \infty$  (22)

**Point distribution /  $\delta$ -function:**  $\delta(x - a) \equiv \lim_{\epsilon \rightarrow 0} \frac{1}{\sqrt{2\pi\epsilon^2}} e^{-\frac{(x-a)^2}{2\epsilon^2}}$  (23)

such that

$$\int_{-\infty}^y f(x)\delta(x - a)dx = \begin{cases} f(a) & \text{if } y > a \\ 0 & \text{if } y < a \end{cases} \quad (24)$$

This allows us to write a discrete distribution in continuous terms:

$$p(x)dx = \sum_k \Pr[X = k] \cdot \delta(x - k)dx \quad (25)$$

- **Joint distributions**,  $p(x, y)$ , including

**Marginalization / law of total probability:**  $p(x) = \int p(x, y)dy$  (26)

**Conditional probability:**  $p(x|y) = \frac{p(x, y)}{p(y)}$  (27)

**Independence:**  $p(x, y) = p(x)p(y)$  (28)

- **Expectation values**, e.g.

**Expected value / ensemble average:**  $\mathbb{E}[g(X)] \equiv \langle g(X) \rangle \equiv \int g(x)p(x) dx$  (29)

**Conditional expectation:**  $\mathbb{E}[g(X)|Y] \equiv \langle g(X) \rangle \equiv \int g(x)p(x|y) dx$  (30)

**Mean:**  $\mu_X \equiv \mathbb{E}[X] \equiv \langle X \rangle \equiv \int xp(x)dx$  (31)

**Variance:**  $\sigma_X^2 \equiv \text{Var}(X) \equiv \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$  (32)

**Moment generating function / Laplace transform:**  $H(z) \equiv \mathbb{E}[e^{-zX}]$  (33)

From the properties of Taylor series, the definition of the moment generating function implies that

$$H(z) = \sum_{k=0}^{\infty} \frac{(-z)^k}{k!} \langle X^k \rangle \quad (34)$$

or

$$\langle X^k \rangle = (-1)^k \left. \frac{\partial^k H}{\partial z^k} \right|_{z=0} \quad (35)$$

- **Central limit theorem:** For independent and identically distributed random variables  $X_1, \dots, X_n$  with a finite mean and variance, define the **sample average**  $\bar{X} \equiv \frac{1}{n} \sum_{i=1}^n X_i$ . Then

$$\mathbb{E}[\bar{X}] = \mu_X, \quad \text{Var}(\bar{X}) = \frac{\sigma_X^2}{n} \quad (36)$$

and, as  $n \rightarrow \infty$ , the distribution of  $\bar{X}$  converges to a Gaussian distribution with mean  $\mu_{\bar{X}} = \mu_X$  and variance  $\sigma_{\bar{X}}^2 = \sigma^2/n$ .

**Note:** A similar approximation applies for any random variable whose distribution is sufficiently “**strongly peaked**” around a characteristic value  $x_0$ . In this case, Taylor expansion of  $\ell(x) \equiv \log p(x)$  yields:

$$\ell(x) \approx \ell(x_0) + \ell'(x_0)(x - x_0) + \frac{1}{2}\ell''(x_0)(x - x_0)^2 = \ell(x_0) - \frac{|\ell''(x_0)|}{2}(x - x_0)^2 \quad (37)$$

so that  $p(x)$  is approximately Gaussian with mean  $\mu_X = x_0$  and variance  $\sigma_X^2 = |\ell''(x_0)|^{-1}$ .

## Programming Prerequisites

The problem sets will also make frequent use of basic programming skills, e.g. the ability to read numerical or textual data from a file, do a few calculations with it, and plot the results in a figure. This is best illustrated with an example from the first problem set. If you feel confident carrying out the programming tasks in this problem in your favorite language (or feel like you could quickly learn these things as you go using the internet), then you will be in good shape for the rest of the course as well.

### Example problem: molecular evolution and genetic diversity in the influenza virus

The text file `influenza_HA_dna_sequences.fasta` contains a list of 841 complete DNA sequences of the hemagglutinin (HA) gene in influenza virus samples collected between 1968 and 2005. Hemagglutinin is a surface protein that allows the viruses to enter host cells, making it a primary target for neutralizing antibodies. This creates a strong selection pressure for the HA gene to evolve over time to evade these immune defenses.

- Calculate the number of single nucleotide differences between the first sample (A/Aichi/2/1968) and the remaining samples, and plot the results as a function of the sampling year. How many differences have accumulated over this  $\sim 40$  year period? What fraction of the HA gene does this account for?
- Calculate the number of genetic differences between all pairs of strains from the same year, and plot the distribution of this quantity aggregated across all years. Estimate the genetic “turnover time” – i.e., how long would we have to wait for the population to accumulate the same number of genetic differences that typically separate co-circulating strains.