

Statistical Bioinformatics, Makerere Basic Probability Timo Koski

TK

02.08.2018



Why ?

Why should there be probabilistic modelling in analysis of biological databases ?

In the literature on data mining and *machine learning* (Baldi and Brunak 1998) we frequently find *Bayesian probability* advocated as the fundamental tool of uncertain reasoning.

- P. Baldi and S. Brunak (1998): *Bioinformatics. The Machine Learning Approach*. A Bradford Book. The MIT Press. Cambridge Massachusetts, London, England.



Why ?

The Cox-Jaynes axioms¹ of uncertain reasoning are often cited in this context, cf. (Arnborg and Sjödin 2000), and would thus seem to provide rationale for the infusion of probability in bioinformatics.

- S. Arnborg and G. Sjödin. (2000): On the Foundations of Bayesianism. *MaxEnt 2000: The Twentieth International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*. American Institute of Physics, pp. 61–71.

¹<http://omega.albany.edu:8008/JaynesBook.html>

The text by Ewens and Grant does not subscribe to much of a Bayesian view of statistics. They do not refer to machine learning or data mining. We are going to learn about, e.g., significance of the BLAST score, about p-values , e.t.c. .

\mathcal{X} and \mathcal{Y} are two discrete *alphabets*, whose generic elements are called *symbols* and denoted by x_i and y_j , respectively.

$$\mathcal{X} = \{x_1, \dots, x_L\}, \mathcal{Y} = \{y_1, \dots, y_J\}.$$

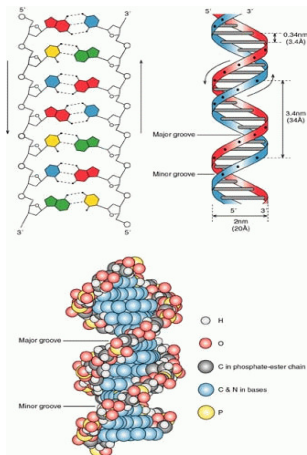
$|\mathcal{X}| \stackrel{\text{def}}{=} \text{the number of elements in } \mathcal{X} = \text{cardinality of } \mathcal{X} = L \leq \infty,$
 $|\mathcal{Y}| = J \leq \infty.$ Unless otherwise stated the alphabets considered here are finite.

In analysis and modelling of DNA sequences the alphabet is

$$\mathcal{X} = \{A, T, C, G\}.$$

These are the four bases present in subunits of DNA called *nucleotides*: adenine (A), thymine (T), cytosine (C) and guanine (G). The particular order of these bases arranged along the sugarphosphate backbone is called the *DNA sequence*. □

DNA



A *sequence* or a *string* \mathbf{x} of length m is an ordered list of m symbols from an alphabet \mathcal{X} written contiguously from left to right

$$\mathbf{x} = x_{l_1} x_{l_2} \dots x_{l_i} \dots x_{l_m}; x_{l_i} \in \mathcal{X}, i = 1, \dots, m.$$

The length of the symbol is denoted by $|\mathbf{x}| = m$. The *empty sequence* \emptyset has length 0 and contains no symbols taken from \mathcal{X} .

Alphabet, Sequence

DNA's genetic code can be represented as a single alphabetic sequence composed of these four symbols. It is by means of this code that the gene controls the formation of other substances in the cell. Progress of biotechnology as well as of computer hardware and software have made it possible to determine the nucleotide sequences for large and increasing number of organisms. Such data are now available in computer-readable form, so it is possible to look for and analyze patterns within sequences using probabilistic modelling and computer algorithms.



There are numerous *model families* assigning a probability $P(\mathbf{x})$ to a string (or by extension to a set of strings) that have been applied in bioinformatics. The model families and the methods of probabilistic modeling often used in bioinformatics are presented in this course.

Random Variables (R.V's)

X is a (**discrete**) random variable that assumes values in \mathcal{X} and Y is a (**discrete**) random variable that assumes values in \mathcal{Y} .



Events are often simply written as

$\{X = x_i\}$: X assumes the value x_i

$\{Y = y_j\}$: Y assumes the value y_j .

Then the probability of the event $\{X = x_i\}$ is

$$f_X(x_i) \stackrel{\text{def}}{=} P(X = x_i).$$

The probability of the event $\{Y = y_j\}$ is

$$f_Y(y_j) \stackrel{\text{def}}{=} P(Y = y_j).$$

Furthermore

$$f_X \stackrel{\text{def}}{=} (f_X(x_1), \dots, f_X(x_L))$$

designates a discrete probability distribution on \mathcal{X} and

$$f_Y \stackrel{\text{def}}{=} (f_Y(y_1), \dots, f_Y(y_J))$$

designates a discrete probability distribution on \mathcal{Y} .

$$f_X(x_i) \geq 0 \quad (1)$$

$$\sum_{i=1}^L f_X(x_i) = 1 \quad (2)$$

and similarly for f_Y .

If $A \subseteq X$ then

$$P_X(A) \stackrel{\text{def}}{=} \sum_{x_i \in A} f_X(x_i) \quad (3)$$

is the probability of the event that X assumes a value in A , a subset of \mathcal{X} .
If $x_1 \leq x_2 \leq \dots \leq x_L$, these are real numbers,

$$F_X(x) = \sum_{i: x_i \leq x} f_X(x_i),$$

is called the distribution function.

If $x_1 \leq x_2 \leq \dots \leq x_L$, real numbers,

$$F_X(x) = \sum_{i: x_i \leq x} f_X(x_i)$$

is called the distribution function.

- $0 \leq F_X(x) \leq 1$.
- $x \leq x', F_X(x) \leq F_X(x')$

$$A \subseteq X$$

$$P_X(A) \stackrel{\text{def}}{=} \sum_{x_i \in A} f_X(x_i)$$

From this one easily finds the complement rule

$$P(A^c) = 1 - P(A), \quad (4)$$

where A^c is the complement of A , i.e. those outcomes that do not lie in A . Also

$$P(A \cup B) = P(A) + P(B) - P(A \cap B), \quad (5)$$

is immediate.

The notation f_X encompasses *tabular* probability distributions, which are simply tables of numbers with the stated properties, and *algorithmic* distributions, which are algorithms for computing probabilities. There are both sorts of probability distributions in bioinformatics as well as combinations thereof.

A *parametric* distribution is a special kind of algorithmic distribution: it consists of a few numerical parameters and a relatively simple algorithm, usually a formula for computing probabilities given some specific values of these parameters. The examples immediately following belong to the parametric type.



Consider the binary alphabet $\mathcal{X} = \{0, 1\}$. Let p be a number between zero and one. If X is a random variable assuming values in $\{0, 1\}$ and if

$$f_X(1) = P(X = 1) = p$$

then we call X a Bernoulli random variable with the *probability of success* p . We write

$$X \in Be(p).$$

We refer to p as a *parameter* of the distribution f_X .

Consider an infinite sequence of Bernoulli random variables $\{X_n\}_{n=1}^{\infty}$ with a common probability of success p . This models an infinite sequence of zeros and ones. Let

X = the number of trials before the first failure not including this.

Then $\mathcal{X} = \{0, 1, 2, 3, \dots, \infty\}$. We say that X has a *geometric probability distribution* if

$$f_X(k) = (1-p) \cdot p^k; k = 0, 1, \dots,$$

This is written as $X \in Ge(p)$.

Suppose $x \in \{0, 1, 2, 3, \dots\}$ and suppose that

$$1 - F_X(x-1) \sim Cp^x$$

for some fixed constant C , $0 < C < 1$. We say that X has a *geometric-like probability distribution*.

$$1 - F_X(x-1) \sim Cp^x$$

We say that the function f is *asymptotic* to function g , as $t \rightarrow \infty$, and write $f \sim g$, if

$$\lim_{t \rightarrow \infty} \frac{f(t)}{g(t)} = 1.$$

Let $\mathcal{X} = \{x_1, x_2, \dots, x_L\}$ be a finite alphabet and

$$f_X(x_k) = \frac{1}{L}; k = 1, 2, \dots, L.$$

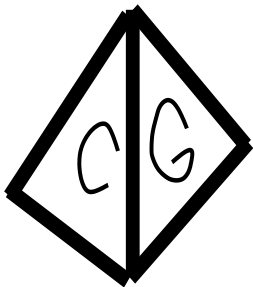
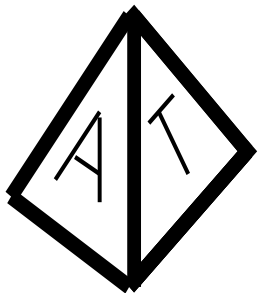
This is called the *uniform probability distribution*.

If X is a random variable assuming values in this alphabet and $P(X = k) = \frac{1}{L}$, then we say that

' X is selected at random'

and we write

$$X \in U(1, L).$$



are not generated like this !

$U(1, 4)$: DNA sequences

Let

$$f_X(k) = e^{-\lambda} \cdot \frac{\lambda^k}{k!}, k = 0, 1, 2, \dots$$

This is a probability distribution on the non-negative integers called the *Poisson distribution*. A random variable X with this distribution is said to be Poisson distributed and we write

$$X \in Po(\lambda).$$

Let A be some event defined in terms of some basic experiment. Introduce the random variable $X =$ *the number of times the event A occurs in n repetitions of some basic experiment*. Then the alphabet of X is $\mathcal{X} = \{0, 1, \dots, n\}$. If the repetitions are modeled as independent (?), then X has the *binomial distribution*:

$$f_X(k) = P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

Binomial Distribution: The Conditions

A **binomial distribution** results from a procedure that meets all the following conditions:

The procedure has a fixed number of random events.

1

The events have outcomes in two categories.

2

The events are independent.

3

The probabilities are constant for each event.

4

$$f_X(k) = P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

(Recall that $\binom{n}{k} = \frac{n!}{(n-k)!k!}$ (the *binomial coefficients*), $0! \stackrel{\text{def}}{=} 1$). We say that X is a binomial random variable with parameters p and n and write

$$X \in \text{Bin}(n, p).$$

Clearly the distribution of X is such that

$$X = X_1 + X_2 + \dots + X_n,$$

as an equality in distribution, where X_1, X_2, \dots, X_n are independent $Be(p)$ -variables as in the preceding example.

The number of matches in a sequence comparison ?

$$X = X_1 + X_2 + \dots + X_n,$$

where X_1, X_2, \dots, X_n are independent $Be(p)$ - variables.

Alignment of Sequences & series of random events

We wish to compare two sequences \mathbf{x} and \mathbf{y} with 15 nucleotides in each.



We say that we have a **match**, if the paired nucleotides are the same in both sequences. We have eleven matches indicated by \downarrow .

Alignment of Sequences & Binomial Distribution

$$\begin{array}{r} \mathbf{x} \\ \mathbf{y} \end{array} = \begin{array}{ccccccccccccccc} & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ G & A & T & A & A & G & C & C & C & C & T & G & T & C & T \\ C & A & A & A & A & T & C & C & C & C & A & G & T & C & T \end{array}$$

If we are willing to assume that the nucleotides are random (DNA die !) and independent, and that the probabilities are the same at each site, then **the probabilities of the number of matches (=successes) follow a binomial distribution with parameters $n = 15$, $p = \frac{1}{4}$** . Then we can compute how probable or likely it is to get 11 matches in a sequence of 15 nucleotides. This is a case of the question of significance (**p-value**) in evaluating sequence alignments.

If, under a given assumption, the probability of an observed event is extremely small (≈ 0), we conclude that the assumption is likely not correct.

In the example with alignment of sequences, the probability of 11 matches in two sequences of 15 nucleotides under the assumption of independent tosses of the DNA die is $1.0297e - 04$. (A computation done using a Matlab function for the binomial probability `binopdf(x,n,p)` with $n = 15$, $p = \frac{1}{4}$, $x = 11$).

If np is small, i.e. if the event A occurs 'rarely', the binomial distribution can be approximated by a Poisson distribution $Po(\lambda)$ with $\lambda = np$ (*the law of small numbers*). The approximation $Bin(n, p) \approx Po(np)$, or more appropriately its generalizations have turned out to be useful for a number of problems in bioinformatics.

Negative Binomial R.V.'s

Same setting as in Binomial R.V.'s but the number of successes is fixed in advance at m and the number of experiments up to and including the m th success is a random variable X .

$$f_X(n) = P(X = n) = \binom{n-1}{m-1} p^m (1-p)^{n-m}.$$

$n = m, m + 1, \dots$. X is a negative binomial random variable with parameters p and m .



The number of failures is fixed in advance at $k + 1$ and the number of experiments up to and not including the $k + 1$ st failure is a random variable X . Then as in the above

$$f_X(x) = P(X = x) = \binom{x}{k} p^{x-k} (1-p)^{k+1}.$$

$x = k, k + 1, \dots$. This reduces to a geometric distribution when $k = 0$, hence this is a generalized geometric distribution.

Hypergeometric R.V.'s

n balls in an urn, r are red and $n - r$ are black. Introduce for $m \leq n$ the random variable $X =$ *the number of red balls is k , if m balls are chosen at random and with replacement, X has the hypergeometric distribution:*

$$P(X = k) = \frac{\binom{r}{k} \binom{n-r}{m-k}}{\binom{n}{m}}.$$

$$X \in \text{Hyp}(n, m, r).$$

R. Román-Roldán, P. Benaola-Galván and J.L. Oliver: Sequence Compositional Complexity of DNA through an entropic segmentation method. *Physical Review Letters*, 80 (6), 1344–1347.

If \mathcal{X} is a discrete alphabet and a finite set of real numbers we define the expectation, $E(X)$, of the random variable X as

$$E(X) \stackrel{\text{def}}{=} \sum_{i=1}^L x_i f_X(x_i).$$

In case $L = \infty$ we require absolute convergence of the series.

If \mathcal{X} is an arbitrary alphabet and $g(\cdot)$ is a map (function) of \mathcal{X} to R ,
 $g : \mathcal{X} \mapsto R$,

$$E(g(X)) \stackrel{\text{def}}{=} \sum_{i=1}^L g(x_i) f_X(x_i).$$

Let \mathcal{X} be a real and finite alphabet and let the *variance*, $\text{Var}(X)$, of X be

$$\begin{aligned}\text{Var}(X) &\stackrel{\text{def}}{=} E([X - E(X)]^2) \\ &= \sum_{i=1}^L (x_i - E(X))^2 f_X(x_i)\end{aligned}\quad (6)$$

Hence $\text{Var}(X) = E(X^2) - (E(X))^2$.

Chebysjev's inequality is another partial intuitive interpretation of variance.
For $k > 0$,

$$P(|X - E(X)| > k) \leq \frac{\text{Var}(X)}{k^2}.$$

Chebysjev's inequality: proof

$$\begin{aligned}\text{Var}(X) &= \sum_{i=1}^L (x_i - E(X))^2 f_X(x_i) = \\ &= \sum_{\{x_i \mid |x_i - E(X)| > k\}} (x_i - E(X))^2 f_X(x_i) + \\ &+ \sum_{\{x_i \mid |x_i - E(X)| \leq k\}} (x_i - E(X))^2 f_X(x_i) \\ &\geq k^2 \sum_{\{x_i \mid |x_i - E(X)| > k\}} f_X(x_i) \\ &= k^2 P(|X - E(X)| > k),\end{aligned}\tag{7}$$

where we have used the fact that $\sum_{\{x_i \mid |x_i - E(X)| \leq k\}} (x_i - E(X))^2 f_X(x_i) \geq 0$. □

X is an integer-valued discrete R.V.,

$$\begin{aligned}\mu_{[r]} &\stackrel{\text{def}}{=} E[X(X-1)\cdots(X-r+1)] = \\ &= \sum_{x:\text{integer}} (x(x-1)\cdots(x-r+1)) f_X(x).\end{aligned}$$

is called the r :th factorial moment.

X is an integer-valued discrete R.V..

$$E \binom{X}{r} = E [X(X-1)\cdots(X-r+1)] / r!$$

is called the binomial moment.

Discrete random variables are in ways more relevant to bioinformatics than continuous random variables. But:

- Continuous variables are needed for approximations.
- Continuous variables are needed in Bayesian statistics.

A distribution function $F_X(x) \stackrel{\text{def}}{=} P(X \leq x)$ such that

$$F_X(x) = \int_{-\infty}^x f_X(t) dt$$

where $f_X(x)$ is a density function,

- $f_X(x) \geq 0$.
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$.

Expectation, $E(X)$, is defined as

$$E(X) = \int_{-\infty}^{+\infty} xf_X(x)dx$$



Variance $Var(X)$ is defined as

$$Var(X) = \int_{-\infty}^{+\infty} (x - E(X))^2 f_X(x) dx$$

$$Var(X) = E(X^2) - (E(X))^2.$$

The random variable $Y = g(X)$ has the expectation

$$E[Y] = E[g(X)] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx.$$

Let $a < b$ and

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{elsewhere.} \end{cases}$$

The $X \in U(a, b)$, uniformly distributed.

If $a < b$ and $X \in U(a, b)$,

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{elsewhere} \end{cases}$$

$$E(X) = \frac{a+b}{2}, V(X) = \frac{(b-a)^2}{12}$$

X has the density $\phi(x; \mu, \sigma)$ defined by

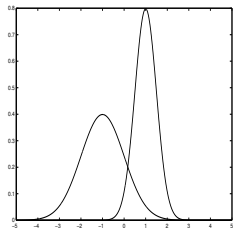
$$\phi(x; \mu, \sigma) \stackrel{\text{def}}{=} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

X has a normal distribution $N(\mu, \sigma^2)$.

Normal random variables

X has a normal distribution $N(\mu, \sigma^2)$.

$$E(X) = \mu, V(X) = \sigma^2.$$



$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & 0 < x < \infty \\ 0 & \text{elsewhere.} \end{cases}$$

$$F_X(x) = 1 - e^{-\lambda x}$$

$X \in \text{Exp}(\lambda)$.

Exponential and Geometric R.V's

$$X \in \text{Exp}(\lambda)$$

$$Y = \lfloor X \rfloor$$

$$k = 0, 1, \dots,$$

$$\begin{aligned} \text{Prob}(Y = k) &= \text{Prob}(k \leq X < k + 1) \\ &= (1 - e^{-\lambda}) e^{-\lambda k}. \end{aligned}$$

Hence $Y \in \text{Ge}(e^{-\lambda})$.



Density of the Fractional Part

$X \in \text{Exp}(\lambda)$

$$D = X - \lfloor X \rfloor$$

D is the fractional part of X .

$$f_D(d) = \begin{cases} \frac{\lambda e^{-\lambda d}}{1 - e^{-\lambda}} & 0 < d < 1 \\ 0 & \text{elsewhere.} \end{cases}$$

$$f_X(x) = \begin{cases} \frac{\lambda^k x^{k-1} e^{-\lambda x}}{\Gamma(k)} & 0 < x < \infty \\ 0 & \text{elsewhere.} \end{cases}$$

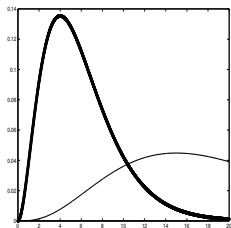
λ and k are positive parameters. $\Gamma(k)$ is Euler's Gamma function. (See Westergren & Råde: **BETA. Mathematics Handbook** for facts and formulae about Euler's Gamma function.)

-

$$\Gamma(z) = \int_1^{\infty} t^{z-1} e^{-t} dt$$

- $\Gamma(z + 1) = z\Gamma(z)$
- $\Gamma(k) = (k - 1)!$, if k is a positive integer.
- $\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$.

Gamma Density



$$f_X(x) = \begin{cases} \frac{\lambda^k x^{k-1} e^{-\lambda x}}{\Gamma(k)} & 0 < x < \infty \\ 0 & \text{elsewhere.} \end{cases}$$

$\lambda = 1/2$, $k = (1/2)\nu$, Gamma density becomes

$$f_X(x) = \begin{cases} \frac{x^{\frac{1}{2}\nu-1} \lambda e^{-\frac{1}{2}x}}{2^{\frac{1}{2}\nu} \Gamma(\frac{1}{2}\nu)} & 0 < x < \infty \\ 0 & \text{elsewhere.} \end{cases}$$

Chi-Square with ν degrees of freedom. This is important in statistical testing.

$$f_X(x) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} x^{\beta-1} & 0 < x < 1 \\ 0 & \text{elsewhere.} \end{cases}$$

$X \in \text{Beta}(\alpha, \beta)$. This density is important in Bayesian statistics.

Probability Generating Function

Let X have values $k = 0, 1, 2, \dots$.

$$G(t) = E\left(t^X\right) = \sum_{k=0}^{\infty} t^k f_X(k)$$

is called the probability generating function.





$$\begin{aligned}\frac{d}{dt} G(t) &= \sum_{k=1}^{\infty} k t^{k-1} f_X(k) \Big|_{t=1} \\ &= E[X]\end{aligned}$$



$$\mu_{[r]} = E[X(X-1)\cdots(X-r+1)] = \frac{d^r}{dt^r} G(t) \Big|_{t=1}$$

$Z = X + Y$, X and Y integer valued, independent (?),



$$G_Z(t) = E(t^Z) =$$
$$E(t^{X+Y}) = E(t^X) \cdot E(t^Y) = G_X(t) \cdot G_Y(t).$$

- $X \in \text{Be}(p)$

$$G_X(t) = 1 - p + pt.$$

- $Y \in \text{Bin}(n, p)$

$$G_Y(t) = (1 - p + pt)^n$$

$$G(t) = E(t^X) = \sum_{k=0}^{\infty} t^k f_X(k)$$

is also known as called the *Mellin transform*.

Moment Generating Function

$$\phi_X(s) \stackrel{\text{def}}{=} E(e^{sX}) = \begin{cases} \sum_{x_i} e^{sx_i} f_X(x_i) & X \text{ discrete} \\ \int_{-\infty}^{\infty} e^{sx} f_X(x) dx & X \text{ continuous} \end{cases}$$

is called the moment generating function.



$$\frac{d}{ds}\phi(0) = E[X]$$



$$\phi(0) = 1$$

$$\frac{d^k}{ds^k}\phi(0) = E[X^k].$$

- $X \in N(\mu, \sigma^2)$

$$\phi_X(s) = e^{\mu s + \frac{1}{2}\sigma^2 s^2}$$

- $Y \in \exp(\lambda)$

$$\phi_Y(s) = \frac{\lambda}{\lambda - s}$$

Let X be a discrete R.V. with moment generating function $\phi_X(s)$. Say that X can take at least one negative value (say $-a$) with positive probability $f_X(-a)$ and at least one positive value (say b) with positive probability $f_X(b)$, and that the expectation of X is nonzero. Then there exists a unique nonzero value s^* such that

$$\phi_X(s^*) = 1.$$

Proof (1)

Proof: Assume $\phi_X(s)$ is defined for all s .

$$\phi_X(s) > f_X(-a)e^{-as}, \phi_X(s) > f_X(b)e^{bs}$$

Hence $\phi_X(s) \rightarrow \infty$, as $s \rightarrow \infty$ and also as $s \rightarrow -\infty$.

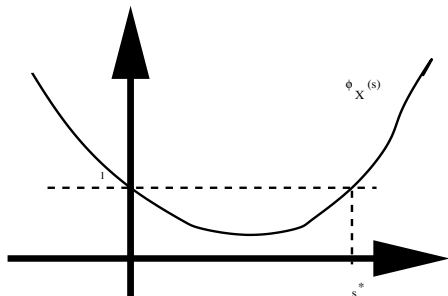


Proof (2)

$$\frac{d^2}{ds^2}\phi(s) = \sum x^2 e^{sx} f_X(x) \geq 0.$$

Hence $\phi(s)$ is convex as function of s . $\phi(0) = 1$ and the the expectation $\frac{d}{ds}\phi(0) = E[X]$ is nonzero by assumption. If $E[X]$ is negative, then the graph of $\phi_X(s)$ must be as in the Figure below:

Proof (3)



Similarly for $E[X]$ positive. □

End of Lecture

