In [1]:	Нуклеотидын IUPAC тэмдэглэгээ Тэмдэглэгээ International Union of Pure and Applied Chemistry-IUPAC 1.2 Basic Algorithms def validate_dna(dna_seq):
	<pre>""" Checks if DNA sequence is valid. Returns True is sequence is valid, or False otherwise. """ seqm = dna_seq.upper() valid = seqm.count("A") + seqm.count("C") + seqm.count("T") if valid == len(seqm): return True else: return False validate_dna("atagagagatctcg") validate_dna("ATAGAXTAGAT")</pre>
Out[1]: In [2]: Out[2]:	<pre>def frequency(seq): """ Calculates the frequency of each symbol in the sequence. Returns a dictionary. """ dic = {} for s in seq.upper(): if s in dic: dic[s] += 1 else: dic[s] = 1 return dic frequency("atagataactcgcatag") #frequency("MVVMKKSHHVLHSQSLIK")</pre>
In [3]:	# Амин хүчлийн давтамжийг тоолно. # Давтамжийг эрэмбээр хэвлэнэ. # lambda функц хэрэглэсэн seq_aa = input("Protein sequence:") freq_aa = frequency(seq_aa) list_f = sorted(freq_aa.items(), key=lambda x: x[1], reverse = True) for(k,v) in list_f: print("Aminoacid:", k, ":", v) Protein sequence:MVVMKKSHHVLHSQSLIK Aminoacid: V: 3 Aminoacid: K: 3 Aminoacid: S: 3 Aminoacid: S: 3 Aminoacid: H: 3
In [4]:	Aminoacid: M: 2 Aminoacid: L: 2 Aminoacid: Q: 1 Aminoacid: I: 1 #Дарааллын урттай харьцуулахад 'G' ба 'C' нуклеотидын хувь def gc_content (dna_seq): """ Returns the percentage of G and C nucleotides in a DNA sequence. """ gc_count = 0 for s in dna_seq: if s in "GCgc": gc_count += 1 return gc_count / len(dna_seq) gc_content("atagataactcgcatag")
Out[4]: In [5]:	0.35294117647058826
In [6]:	2. Хөрвүүлэлт (Transcription) ба урвуу гүйцээлт(Transcription and Reverse Complement) днх-ийг рнх молекул руу хөрвүүлэх: днх-ийн молекулуудад агуулагдах генетик мэдээлэл нь уургийн нийлэгжих бүх процессын үндэс болдог.
111 [0].	# DNA-ийн молекулууд нь урвуу чиглэлд уншигдах хоёр цуваатай. # Хөрвүүлэлтийг үүсэх үед уг хоёр цуваа салж, шинэ RNA молекул нь нэг цувааны гүйцээлт болж үүснэ. # Тиймээс ДНХ-ийн нөгөө цуваатай төстэй (бид үүнийг template цуваа гэдэг), 'Т'-ийн оронд 'U' байна. def transcription (dna_seq): """ Function that computes the RNA corresponding to the transcription of the DNA sequence provided. """ assert validate_dna(dna_seq), "Invalid DNA sequence" return dna_seq.upper().replace("T","U") transcription("atagataactcgcatag")
Out[6]:	
	elif c== "C":
In [8]:	Messenger RNA молекулуудад агуулагдах мэдээллийн дагуу амин хүчлүүдийн гинжийн дарааллыг үүсгэн уураг нүйгээ. # Самнех алгоритмаас харахад DNA-ээс mRNA үүсгэх нь хялбар тул DNA-ээс шууд уураг үүсгэе. # Тулхүүрүүд нь кодонууд (нийт 64), холбогдох утгууд нь амин хүчлүүд. # Буруу кодоныг аргумент болгон дамкуулсан тохиолдолд функц нь None-г буцаана. # Зогсох кодомд '_ тэмдгийг ашигласан def translates a codon into an aminoacid using an internal dictionary with the standard genetic code.""" tc = ("GCT":"A", "GCC":"A", "GCA":"A", "GCG":"A", "GCG":"A", "GCG":"C", "GAT":"C", "GCG":"C", "GGG":"C", "GG
In [9]:	else: return None # ДНХ-ийн дараалал болон хувиргах анхны байрлалыг оролт болгон авч, холбогдох амин хүчлийн дарааллыг буцаана. # Дарааллыг кодон болгон хуваах нь хамгийн сүүлийн аргумент нь алхамын утгыг тодорхойлох боломжийг олгодог # гапде функцийг ашиглан тохирох индексүүдийг үүсгэх замаар хийгдэнэ. # Хуваасан кодонуудын хувиргалт нь өмнө боловсруулсан функцийг буцаах утгаар ирнэ def translate_seq (dna_seq, ini_pos = 0): """ Translates a DNA sequence into an aminoacid sequence. """ assert validate_dna(dna_seq), "Invalid DNA sequence" seqm = dna_seq.upper() seq_aa = "" for pos in range(ini_pos,len(seqm)-2,3): cod = seqm[pos:pos+3] seq_aa += translate_codon(cod) return seq_aa
In [10]:	Print(translate_seq('ATAGATACTCGCATAG')) Print(translate_seq('ATAGATACTCGCATAG')) Freether K KOA HS ONDH CONFONTON TYN ABSTARACH YTTYYA GAÄHA, Whis: HOY AMMH XYYNHÄF BEQ BEQS KOAQHOLOD KOANDADOF: YYNHÄF KOAQOHM XODOFNOO (CODON KOADADADOF: # ONDHAF HORDH KOAQOHM XODOFNOO (CODON KOADADADOF: # ONDHAF HORDH KOAQOHM XODOFNOO (CODON HORDH KOAQOHM) YTA HS KOAQOH KOYDHÄH XODOFNOO (MARINE MARINE) WHO Provides the frequency of each codon encoding a given aminoacid, in a DNA sequence .""" assert validate_dna(dna_seq), "Invalid DNA sequence" seqm = dna_seq.upper() dic = {} total = 0 for in range(0, len(seqm)-2, 3): cod = seqm[i:i+3] if translate_codon(cod) == aa: if cod in dic: dic[cod] += 1 else: dic[cod] = 1 total += 1 if total >0: for k in dic: dic[k] /= total return dic
In [11]:	 4. Таамаг ген (Putative Genes) илрүүлэх Ерөнхий тохиолдолд DNA дарааллын (жнь, genome sequencing project) кодлох бүсүүд хаана байгаа нь урьдчилан мэддэгддэггүй. Уургийн хувиргалт нь үргэлж Methionine ("М") амин хүчлийн кодоноор (эхлэх кодон - "ATG") эхэлдэг. Энэ амин хүчил нь уургийн эхлэлээс гадна бусад байрлалд ч тохиолдож болно. Мөн зогсох кодон олдвол хувиргалты процесс дуусдаг байх ёстой. Уг сэжигтэй хэсгүүдийг хайхдаа эхлээд DNA (эсвэл RNA) дараалал дээрх унших хүрээг (reading frames) тооцоолох юм. Унших хүрээ нь ДНХ-ийн дарааллыг дараалсан давхцаагүй гурвалсан (боломжтой кодон) болгон хуваах арга юм.
	# Өгегдсен дараалал нь эхний, хоёр, гурав дахь байрлалаас эхлэн унших боломжтой гурван тохиолдолтой. # Энэ гурваас гадна урвуу гүйцээлтийнх нь гурван хүрээг тооцоолох хэрэгтэй. # Ингэснээр геномын хэсэг дээрх уургийн байж болох бүх кодыг хайх боломжийг олгодог. def reading_frames (dna_seq): """Computes the six reading frames of a DNA sequence (including the reverse complement.""" assert validate_dna(dna_seq), "Invalid DNA sequence" res = [] res.append(translate_seq(dna_seq,0)) res.append(translate_seq(dna_seq,2)) rc = reverse_complement(dna_seq) res.append(translate_seq(c,0)) res.append(translate_seq(rc,1)) res.append(translate_seq(rc,2)) return res reading_frames('ATAGATAACTCGCATAG')
Out[11]:	['IDNSH', '_ITRI', 'R_LA_', 'LCELS', 'YASYL', 'MRVIY'] 4.1 Уншилтын нээлттэй хүрээ хайх Open Reading Frames-ORF хайх: Унших хүрээнүүдийг тооцоолсны дараагийн алхам бол эдгээр хүрээгээр кодлох боломжтой уураг олох явдал юм.
In [12]:	# Амин хүчлийн дарааллаас бүх боломжит уургийг гаргаж авна. # Энэ функц нь дарааллаар явагддаг бөгөөд # 'M' илэрсэн үед таамагласан уургийг current_prot жагсаалтад үүсгэж эхэлнэ # Зогсох тэмдэг олдвол энэ жагсаалтын бүх уураг гаралтын хувьсагчид нэмнэ def all_proteins_rf (aa_seq): """Computes all posible proteins in an aminoacid sequence.""" aa_seq = aa_seq.upper() current_prot = [] proteins = [] for aa in aa_seq: if aa == "_": if current_prot:
	 4.2 Бүх боломжит уургийг тооцоолох Эхлээд унших хүрээнүүдийн хөрвүүлэлтийг тооцоолно Дараа нь зургаан амин хүчлийн дараалал дээр давталт хийнэ Өмнөх функцийг ашиглан бүх боломжит уургийг олж, үр дүнгийн жагсаалтад цуглуулна. Уг функцийн буцаах жагсаалт нь бодит тохиолдолд олон тооны уураг агуулж болзошгүй Жагсаалтыг таамагласан уургийн хэмжээгээр эрэмбэлж, хамгийн бага хэмжээтэйгээс нь эхлэх хэрэгтэй. Жижиг таамаг уургууд нь санамсаргүй байдлаар олдож болох ч хэдэн арван амин хүчил агуулсан уураг ингэж тохиолдох магадлал багатай юм.
In [13]:	• Зогсоох кодон үүсэх магадлал 3/64 буюу 5% орчим байдаг тул ойролцоогоор 20 амин хүчил тутамд зогсолт кодон үүсэх төлөвтэй байдаг
In [14]:	<pre>def all_orfs_ord (dna_seq, minsize = 0): """Computes all possible proteins for all open reading frames. Returns ordered list of proteins with minimum size.""" assert validate_dna(dna_seq), "Invalid DNA sequence" rfs = reading_frames (dna_seq) res = [] for rf in rfs: prots = all_proteins_rf(rf) for p in prots: if len(p) > minsize: insert_prot_ord(p, res) return res def insert_prot_ord (prot, list_prots): i = 0 while i < len(list_prots) and len(prot) < len(list_prots[i]): i += 1 list_prots.insert(i, prot)</pre>
In [15]:	4.3 Putting It All Together
In [16]:	GC content (global): 0.25 Direct translation: M_E All proteins in ORFs (decreasing size): ['M'] 5. БИОЛОГИЙН ДАРААЛЛЫН КЛАСС Илүү зохион байгуулалттай,цаашид хэрэглэх боломжтой КЛАСС бүтэц Объект хандалгат программчлал
	<pre>class MySeq: """ Class for biological sequences. """ definit (self, seq, seq_type = "DNA"): self.seq = seq.upper() self.seq_type = seq_type deflen(self): return len(self, seq) defgetitem(self, n): return self.seq[n] defgetslice(self, i, j): return self.seq[i:j] defstr(self): return self.seq def get_seq_biotype (self): return self.seq_type def show_info_seq (self): print ("Sequence: " + self.seq + " biotype: " + self.seq_type)</pre>
	<pre>def alphabet (self): if (self.seq_type=="DNA"): return "ACGT" elif (self.seq_type=="RNA"): return "ACGU" elif (self.seq_type=="PROTEIN"): return "ACDEFGHIKLMNPQRSTVWY" else: return None def validate (self): alp = self.alphabet() res = True i = 0 while i < len(self.seq) and res: if self.seq[i] not in alp: res = False else: i += 1 return res def transcription (self):</pre>
	<pre>if (self.seq_type == "DNA"):</pre>
	<pre>if (self.seq_type != "DNA"): return None seq_aa = "" for pos in range(iniPos,len(self.seq)-2,3): cod = self.seq[pos:pos+3] seq_aa += translate_codon(cod) return MySeq(seq_aa, "PROTEIN") ## main s1 = MySeq("ATGTGATAAGAATAGAATGCTGAATAAATAGAATGACAT") s2 = MySeq("MKVVLSVQERSVVSLL", "PROTEIN") print(s1.validate(), s2.validate()) print(s1) s3 = s1.transcription() s3.show_info_seq()</pre>
	s4 = s1.reverse_comp().translate() s4.show_info_seq() True True ATGTGATAAGAATAGAATGCTGAATAAATAGAATGACAT Sequence: AUGUGAUAAGAAUAGAAUAGAAUAGAAUAGAAUGACAU biotype: RNA Sequence: MSFYLFSILFLSH biotype: PROTEIN 6. BioPython ашиглан дараалал боловсруулах biopython 6.1 Дарааллын аннотоц объектууд
	Віо. SeqRecord нь дараалал ба тэдгээрийн аннотацын үндсэн контейнер класс. SeqRecord • .seq – дараалал өөрөө буюу Seq классын объект • .id – дарааллын дугаар • .name – дарааллын нэр • .description – дарааллын тайлбар • .annotations – бүхэл дарааллын глобал тайлбарууд, dictionary төрөлтэй, • Түлхүүр: аннотацын төрлүүд болон бүтэцлэгдээгүй талбарууд • Утга: тухайн дарааллын талбаруудын утгууд
	 .features – бүтэцлэгдсэн үзүүлэлтүүд, дараалалд эсвэл хэсгүүдэд нь хамаарах SeqFeature объектын жагсаалт; .letter_annotations – дараалал дахь үсэг (байрлал) бүрийн боломжит аннотацууд; .dbxrefs – өгөгдлийн сангийн лавлах мэдээлэл SeqFeature .location – энэ аннотацад хамаарах дарааллын бүс, FeatureLocation классын объект; .type – тухайн үзүүлэлтийн төрлийн төлөв, тэмдэгт мөр;
In [17]:	 . qualifiers – тухайн үзүүлэлтийн талаар нэмэлт мэдээлэл, талбар-утга dictionary төрөлтэй. 7. Exercises DNA_Codons = {
	"GAT": "D", "GAC": "D", "GAA": "E", "GAGC": "E", "TTT": "F", "TTC": "F", "GGT": "G", "GGG": "G", "GGA": "G", "GGG": "G", "CAT": "H", "CAC": "H", "ATA": "I", "ATT": "I", "ATC": "I", "AAA": "K", "AAG": "K", "TTA": "L", "TTG": "L", "CTT": "L", "CTG": "L", "CTG": "L", "ATG": "M", "AAT": "N", "AAC": "N", "CCT": "P", "CCC": "P", "CCG": "P", "CAAT: "Q", "CAGC": "Q", "CGT": "R", "CGC": "R", "CGA": "R", "AGG": "R", "AGG": "S", "ACT": "S", "TCC": "S", "TCA": "S", "AGG": "S", "AGG": "S", "ACT": "T", "ACC": "T", "ACG": "T", "GTT": "V", "GTG": "V", "GTG": "V", "TGG": "W",
	"TAT": "Y", "TAC": "Y", "TAA": "_", "TAG": "_" }

Биологийн дарааллын үндсэн боловсруулалт(Basic Processing of Biological Sequences)

Олон янзын практик зорилгоор биоинформатикийн алгоритм, хэрэгсэлүүдэд генетикийн мэдээлэл бүхий молекулуудыг нуклеотидын "нэг хэмжээст дараалал" хэлбэрээр

1. Биологийн дараалал: Тэмдэг тэмдэглэгээ, үндсэн алгоритмууд

Symbol

C

D

Ε

F

G

Н

M

Ν

Q

R

S

Т

٧

Т

Nucleotides represented

Name Alanine

Cysteine

Glycine

Histidine

Isoleucine

Methionine

Asparagine

Glutamine

Arginine

Threonine

Tryptophan

Tyrosine

Serine

Valine

Амин хүчлийн IUPAC

Lysine

Proline

Aspartic acid

Glutamic acid

Phenylalanine

1.1 Тэмдэг тэмдэглэгээ

Name

Adenine

Cytosine

Guanine

Thymine

Uracil

Keto

Amino

Purine

Strong

Weak

Not A

Not C

Not G

Not T

Any base

Pyrimidine

C

G

Т

G, T

A, C

A, G

C, G

A, T

C, T

C, G, T

A, G, T

A, C, T

A, C, G

A, C, G, T

дүрсэлдэг.

Symbol

C

G

Т

U

K

M

R

S

W

В

D

Н

else:

return False

for i in self.seq:

pyrim += 1

def codons(self): codon_lst = []

return codon_lst

def isPalindrome(self): revSeq = self.seq[::-1] if(self.seq == revSeq): return "Палиндром мөн" else: return "Палиндром биш"

def CG_content(self):

def checkProtein(self):

for i in range(3):

if(i in DNA_Codons): count += 1

print("Total codons", seq)

if(self.checkProtein != -1):

for i in range(0, len(seq), 3): codon_lst.append(seq[i:i+3])

for i in seq:

if(count > 0): return count else: return -1

def boolProtein(self):

def mappingDict(self): codon_lst = []

for s in codon_lst: if(s in dct): dct[s] += 1 else: dct[s] = 1

def allProteins(self):

seq = self.codons()

aa.append(DNA_Codons[s])

def $k_{req}(self, seq, k = 4)$:

for i in range(0, len(seq), k): codon_lst.append(seq[i:i+k])

sqq = "ACGGACGGACGGAAAAACGGACGGACGGAAAA"

print("CG: " + str(dna.CG_content()) + "%")

print("Уураг мөн үү! " + str(dna.boolProtein())) print("dictionary бүтцээр ", dna.mappingDict()) print("Reading Frame", dna.allProteins())

dictionary бүтцээр {'ACG': 2, 'GTG': 2, 'CTA': 2} Reading Frame ['T', 'V', 'L', 'T', 'V', 'L']

function should return −1 if no protein is found.

К хэмжээгээр амин хүчил гаргах {'ACGG': 6, 'AAAA': 2}

first position of the DNA sequence). Stop codons should be ignored.

2. Write a program that reads a DNA sequence and checks if it is equal to its reverse com- plement.

3. Write and test a function that, given a DNA sequence, returns the total number of "CG" duplets contained in it.

print("Уураг үүсгэж байгаа тоо " + str(dna.checkProtein()))

print("К хэмжээгээр амин хүчил гаргах ", dna.k_freq(sqq, 4))

for s in seq:

codon_lst = []

for s in codon_lst: if(s in dct): dct[s] += 1 else: dct[s] = 1

seq = "ACGGTGCTAACGGTGCTA"

print(dna.isPalindrome())

print(dna.purines_pyrimidines())

return aa

dct = {}

return dct

dna = Bio_seq(seq)

Purines: 10 Pyrim: 8

Уураг мөн үү! True

Уураг үүсгэж байгаа тоо 16

Палиндром биш

CG: 56%

dct = {}

return dct

aa = []

return True else: return False

seq **=** [] count = 0

purines += 1

Кодонуудыг гаргаж авна.

purines = 0pyrim = 0

def purines_pyrimidines(self):

if(i == 'A' or i == 'G'):

elif(i == 'U' or i == 'T' or i == 'C'):

for i in range(0, len(self.seq), 3): codon_lst.append(self.seq[i:i+3])

Дарааллийг палиндром эсэхийг шалгана.

Дарааллын С, G нуклеотидын хувь

return "Purines: " + str(purines) + " Pyrim: " + str(pyrim)

return round((self.seq.count('C') + self.seq.count('G')) / len(self.seq) * 100)

Олдсон дараалал нь уураг мөн үү биш үү гэдгийг логик оператороор тодорхойлдог функц

ДНХ дарааллыг бүрдүүлж буй амин хүчлүүдийн дэд дарааллуудыг олдог функц.

к хэмжээтэй дэд дараалал өгөгдсөн дараалалд оршиж байна уу эсэхийг шалгана.

exam: ACG: T, GTG: V, CTA: L, ACG: T, GTG: V, CTA: L = Thr, Val, Leu, Thr, Val, Leu

Хэрвээ оршиж байвал түлхүүр нь тухайн дэд дараалал ба утга нь хэдэн удаа дэд дараалал давтагдсан тоо болох dictionary буцаана.

Total codons ['ACG', 'GTG', 'CTA', 'ACG', 'GTG', 'CTA', 'CGG', 'TGC', 'TAA', 'CGG', 'TGC', 'GGT', 'GCT', 'AAC', 'GGT', 'GCT']

4. Write and test a Python function that, given a DNA sequence, returns the size of the first protein that can be encoded by that sequence (in any of the three reading frames). The

7. Write a program that reads an aminoacid sequence and a DNA sequence and prints the list of all sub-sequences of the DNA sequence that encode the given protein sequence.

6. Write and test a Python function that, given a DNA sequence, creates a map (dictionary) with the frequencies of the aminoacids it encodes (assuming the translation is initiated in the

8. Write a function that, given a sequence as an argument, allows to detect if there are re-peated sub-sequences of size k (the second argument of the function). The result should be a dictionary where keys are sub-sequences and values are the number of times they oc-cur (at least 2). Use the function in a program that reads the sequence and k and prints the result

1. Write a program that reads a DNA sequence, converts it to capital letters, and counts how many nucleotides are purines and pyrimidines.

5. Write and test a function that given an aminoacid sequence returns a logic value indicat-ing if the sequence can be a protein or not.

Өгөгдсөн амин хүчлийн дарааллууд нь уураг үүсгэж байгаа эсэхийг шалгадаг функц

Эхний давхар давталтаар 3 аар тасдаад 1 алхам хойшлоод 3 аар аваад явна

2 дохь давталтаар кодонуудыг байдаг эсэхийг шалгаж тоолно

Дарааллыг амин хүчлүүдийг dictionary бүтэцээр буцаана

хэрвээ байдаг бол тоог нь үгүй бол -1 буцаана

for j in range(i, len(self.seq), 3): if(len(self.seq[j:j + 3]) == 3):seq.append(self.seq[j:j + 3])

ДНХ-ийн дарааллыг уншиж, том үсэг болгон хувиргаж, хэдэн нуклеотид нь пурин, пиримидин болохыг тоолдог функц