# HYPERDIMENSIONAL COMPUTING FOR PROTEIN LANGUAGE MODELING

...

Michael Fatjanov

Student ID: ...

Supervisor(s): Prof. Dr. Bernard De Baets

GHENT UNIVERSITY

Gent, FILL IN THE DATE

The promotor,                                    The author,

Prof. Dr. Bernard De Baets              Michael Fatjanov

# ACKNOWLEDGEMENTS

Thank you, all of you!

# CONTENTS

# SAMENVATTING

nederlandse samenvatting

# SUMMARY

insert english summary here...

# 1. [INTRODUCTION]

## 1.1 Digital biology, protein sequence research and traditional bioninformatics tools

## 1.2 State-of-the-art, deep learning and protein language modeling

## 1.3 Hyperdimensional computing

Hyperdimensional computing is a relatively new paradigm of computing developed by Kanerva.Kanerva (2009) that tries to mimick the workings of a (human) brain by computing with vectors of tens of thousands elements long, so in the realm of hyperdimensionality. The human brain consists of about 100 billion neurons (nerve cells) and 1000 trillion synapses that connect these neurons. Each neuron is connected up to 10000 other neurons, creating massive circuits. This is likely fundamental to the workings of the human brain and what seperates our brains from modern von Neumann computer architectures which operate on 8 to 64 bit vectors. This becomes clear when we compare the relative simplicity for a human to learn a language compared to computers which require a large and complecated set of arithmetic operations in the form of deep learning networks together with terabytes of data and thousands of Watts of compute power to try to come close mastering a language whilst a human can recognize other languages relatively easy when they don't even speak it. Likewise languages, we can very easily memorize and compare other intrisically complex and contextual concepts such as images. A computer would have a hard time finding similarity between a set of images and faces because this requires very complex

machine learning models. The human brain can do this all with a huge efficiency by uconsuming only roughly 20 W of energy.

To achieve this kind of flexible intelligence, we might have to step away from the restrictive von Neumann architecture and so Kanerva proposes the use of hyperdimensional vectors, a different form of representation for entities by which a computer can compute with. The use of models based on high dimensionality is not entirely new and is being explored since the 1990s. Some of these earlier models include Holographic Reduced RepresentationsPlate (1995), Spatter CodeKanerva (1994) etc. An HDV can represent anything from a scalar number to any kind of concept. This vector is initially made up of totally random elements, but with a simple set of operations which will be explained later, we can use other vectors to combine some concepts into new similar or dissimilar concepts. For example to show the essence of HDC and how it tries to simulate the brain, we can compare the concept of a *table* to the concept of a *brocolli*. We would not immediately conclude that they are in any way similar but as humans we can trace back *table* to *plate* which has some similarities with *food* from which we can easily extract the concept of *brocolli*. These kind of operations are not very obvious for a classical computer but easy for us humans.

The elements in an HDV can be made up from binary bits like in classical computing but also of bipolar or real numbers. The choice of the nature of the elements has also implications on the nature of the different operations and on the results. Highly efficient bit operations could be used on binary vectors but then the amount of information stored in such a vector would be drastically lessened compared to bipolar or real vectors, leading to a lower accuracy.

An initial HDV is made up totally random. This *holistic* or *holographic* representation of a concept smeared out over a vector consisting of thousands of bits gives rise to interesting properties of this paradigm such as its robustness. These kind of systems are very tolerant to noise and failure of bits since we introduce a lot of redundancy in the vector just by stochastics. This is very unlike classical computing where every bit counts and one failure in a bit can lead to disasters. Again, this

# BIBLIOGRAPHY

Kanerva, P. (1994). The spatter code for encoding concepts at many levels. pages 226–229.

Kanerva, P. (2009). Hyperdimensional Computing: An Introduction to Computing in Distributed Representation with High-Dimensional Random Vectors. *Cognitive Computation*, 1(2):139–159.

Plate, T. (1995). Holographic reduced representations. *IEEE Transactions on Neural Networks*, 6(3):623–641.

# APPENDIX A
# [MODEL DEVELOPMENT CODING](#)

## A.1   Pseudocode of the presented algorithm

---

**Algorithm 1:** How to write algorithms

---

**Data:** this text

**Result:** how to write algorithm with LaTeX2e

initialization;

**while** *not at end of this document* **do**

    read current;

    **if** *understand* **then**

        go to next section;

        current section becomes this one;

    **else**

        go back to the beginning of current section;

    **end**

**end**

---

## A.2   Sensitivity base class code

```python
1  import os
2  import numpy as np
3
4  from parameter import *
5  import matplotlib.pyplot as plt
6  from matplotlib.ticker import FixedLocator, MaxNLocator
7
8  class SensitivityAnalysis(object):
9      """
10     Base class for the Sensitivity Analysis
11
12     Parameters
```

```
13          ----------
14      ParsIn : list
15          ModPar class instances in list or list of (min,max,'name')-
                tuples
16
17      Attributes
18          ----------
19      ParsIn : list
20          a list of (min,max,'name') values,
21          [(min,max,'name'),(min,max,'name'),...(min,max,'name')]
22      parmap : dict
23          tracks the sequence of the parameters
24      Pars : list of ModPar instances
25          Used when working with the pyFUSE package
26      ndim :  int
27          number of uncertain input factors
28      namelist : list
29          list of the uncertain input factors used
30
31      """
32
33      def __init__(self,ParsIn):
34          '''
35          Check if all uniform distribution => TODO ! if all -> sobol
                sampling
36          is possible, else, only uniform and normal distribution are
                supported
37          for using the sobol sampling... Here is still work to do!!
38          '''
39
40          if  isinstance(ParsIn, dict): #bridge with pyFUSE!
41              dictlist = []
42              for key, value in ParsIn.iteritems():
43                  dictlist.append(value)
44              ParsIn = dictlist
45              print ParsIn
46
47          #control for other
48          self.ParsIn = ParsIn
49          self.parmap={} #dictionary linking ID and name, since dict
                instance has no intrinsic sequence
50          for i in range(len(ParsIn)):
51              if isinstance(ParsIn[i], ModPar): #or isinstance(ParsIn[i],
                    pyFUSE.parameter.ModPar):
52                  cname = ParsIn[i].name
53                  self.Pars = ParsIn
54                  self.ParsIn[i] = (ParsIn[i].min, ParsIn[i].max, cname)
```

6

```python
55              self.parmap[i] = cname
56
57          elif isinstance(ParsIn[i],tuple):
58              if ParsIn[i][0] > ParsIn[i][1]:
59                  raise Exception('Min value larger than max value')
60              if not isinstance(ParsIn[i][0],float) and isinstance(
                    ParsIn[i][1],float):
61                  raise Exception('Min and Max value need to be float'
                        )
62              if not isinstance(ParsIn[i][2],str):
63                  raise Exception('Name of par needs to be string')
64              self.parmap[i] = ParsIn[i][2]
65              #create modpar instance of the tuple
66              self.Pars=[]
67              for par in ParsIn:
68                  self.Pars.append(ModPar(par[2],par[0],par[1],(par
                        [0]+par[1])/2.,'randomUniform'))
69          else:
70              raise Exception('The input type for sampling not correct
                    ,\
71          choose ModPar instance or list of (min,max)-tuples')
72
73      self.ndim=len(ParsIn)
74
75      self.namelist = []
76      for i in range(self.ndim):
77          self.namelist.append(self.parmap[i])
78
79  def WritePre(self,filename = 'inputparameterfile', *args, **kwargs):
80      '''
81      Parameterinputfile for external model, parameters in the columns
              files
82      and every line the input parameters
83
84      Parameters
85      -----------
86      filename : str
87          name of the textfile to save
88      *args, **kwargs : args
89          arguments passed to the numpy savetext-function
90      '''
91
92      np.savetxt(filename,self.parset2run,*args,**kwargs)
93      print 'file saved in directory %s'%os.getcwd()
94
95  def ReadRuns(self,filename, *args, **kwargs):
96      '''
```

```
97          Read model outputs (TODO: do sobol for multiple outputs,
                iterating the
98          post)
99          Format is: every output of the ithe MC on ith line
100
101         output2evaluate can also be made on a other way
102
103         Parameters
104         -----------
105         filename : str
106             name of the textfile to load
107         *args, **kwargs : args
108             arguments passed to the numpy loadtext-function
109
110         '''
111         self.output2evaluate = np.loadtxt(filename, *args, **kwargs)
```

## A.3   Model input file for PyFUSE model

```
1  #####################################################
2  ##     Model Parameter input file
3  ##     The parameter is defined by his distribution,
4  ##     boundaries and extra info needed by distribution
5  ##     provide on each line one parameter with
6  ##     following information:
7  ##     name : string
8  ##         Name of the parameter
9  ##     minval : float
10 ##         Minimum value of the parameter distribution
11 ##     maxval :  float
12 ##         Maximum value of the parameter distribution
13 ##     optguess : float
14 ##         Optimal guess of the parameter, must be
15 ##         between min and max value
16 ##     pardistribution : string
17 ##         choose a distributionfrom: randomUniform,
18 ##         randomTriangular, randomTrapezoidal,
19 ##         randomNormal, randomLogNormal
20 ##     *kargs  :
21 ##         Extra arguments necessary for the
22 ##         chosen distribution
23 ##   Lines with ## marks are neglected
24 #####################################################
25 ## NAME MIN MAX OPTGUESS DISTRIBUTION ARGS*
26 S1max 50. 5000.000 400. randomTriangular 1000.
```

```
27 S2max 100. 10000.000 1000. randomNormal 500. 25.
28 fitens 0.01 1.0 0.99 randomLogNormal  0.5 0.2
29 firchr 0.050 0.950 0.5 randomTrapezoidal 0.4 0.6
30 fibase 0.050 0.950 0.5 randomUniform
31 r1 0.050 0.950 0.5 randomUniform
32 ku 0.01 1000. 0.044 randomUniform
33 c 0.99 20.0 1. randomUniform
34 alfa 1.000 250. 150. randomUniform
35 psi 1.000 5.0 2.5 randomUniform
36 kappa 0.050 0.950 0.5 randomUniform
37 ki 0.001 1000. 0.00833 randomUniform
38 ks 0.001 10000. 0.5 randomUniform
39 n 1.000 10. 3. randomUniform
40 v 0.00001 0.250 0.004 randomUniform
41 vA 0.001 0.250  0.0015 randomUniform
42 vB 0.001 0.250 0.0015 randomUniform
43 Acmax 0.050 0.950 0.5 randomUniform
44 b 0.001 3.0 0.2 randomUniform
45 loglambda 5.000 10.0 7.5 randomUniform
46 chi 2.000 5.0 3.5 randomUniform
47 mut 0.010 5.0 0.6 randomUniform
48 be 0.99 4. 3.1 randomUniform
49 alfah 0.01 0.99 0.5 randomUniform
50 tg 0.0 0.7 0.3 randomUniform
51 tif 0.0 0.7 0.26 randomUniform
52 tof 0.0 0.7 0.12 randomUniform
53 ko 0.01 0.99 0.15 randomUniform
54 timeo 2. 48 24 randomUniform
55 timei 2. 250. 20 randomUniform
56 timeb 200. 10000. 2100. randomUniform
```