

DATA-PORTFOLIO



Ruperto P. Bonet

ruperto.bonet@gmail.com
Telf:+44(79) 0169-3987

PhD in Computational Mechanics

(CIMEC-UNL, ARGENTINA-1998)

**4.94/5 Bachelor Degree in
Mathematics-Statistics**

(Havana University, Cuba, 1986)

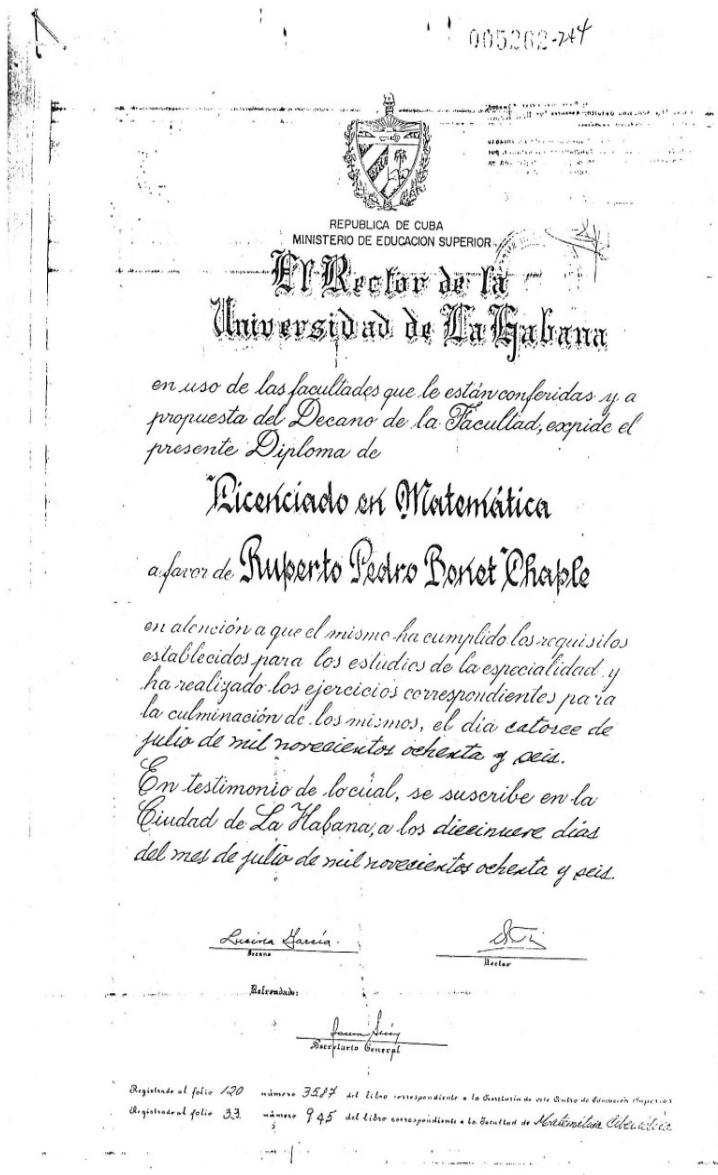
OUTLINE

Personal Description	2
Academic Degrees	3
Research Experiences Fellowship Programs	4
Research Experiences Research Stages	5
Research Experiences: Simulation Projects of engineering problems	6
Publishings	10
Experiences in Statistics	11
AI and Big Data	15
Cloud and Computing	16
Java-Hadoop	17
Spark	18
Experiences as a Big Data Analytics Tutor	19
Machine learning: Classification Training	23
Machine learning: Binary Classification	29
Lumped Sample Classification Prediction Method	34
Machine learning Algorithms	41
Applied AI and Deep Learning Serialization Keras	48
Tourism Enhancers Index by neighborhoods in London Project	53
Applied AI and Image Processing	60
Applied AI and Natural Language Processing	70

PERSONAL DESCRIPTION

- Enthusiastic and successful communicator teacher at all educational levels.
- Passionate scientific researcher in Applied Mathematics focused on Numerical Methods, Computational Mechanics, Simulation and Mathematical Modelling of real-world problems.
- Used to learning new technologies by self and willing to work in different environments and share as well its knowledge.
- A friendly person that participates in very active networking through Facebook, Linkedin, and other social media.
- A person with a cultural wealth acquired living and working around the world.

ACADEMIC DEGREES



Facultad de Ingeniería y Ciencias Hídricas

Por cuanto: el Licenciado en Matemática **Ruberto Pedro Bonet Chaple**,
Pasaporte B-0273184, nacido en La Habana, República de Cuba, el 27 de marzo de 1957, ha
aprobado la carrera de Doctorado en Ingeniería Mención: Mecánica Computacional, el 4 de diciembre
de 1998.

Por tanto: de conformidad con lo que dispone el Estatuto vigente, se le otorga el presente diploma de
Doctor en Ingeniería Mención: Mecánica Computacional.

Santa Fe, 17 de junio de 1999.

M. P. B.
CHRISTIAN V. COELLO
NOTARIO

M. C. V.
MARIA CRISTINA VILLALBA
SECRETARIA ADMINISTRATIVA OFICINA DE LA FACULTAD

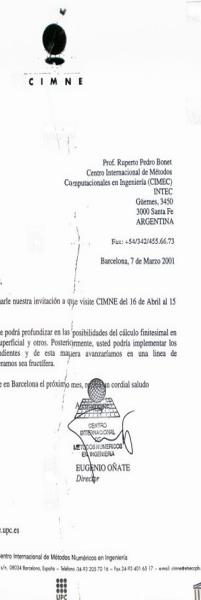
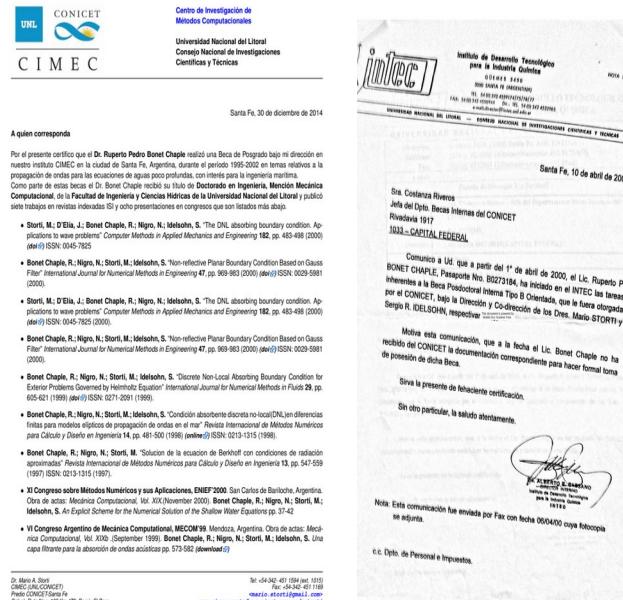
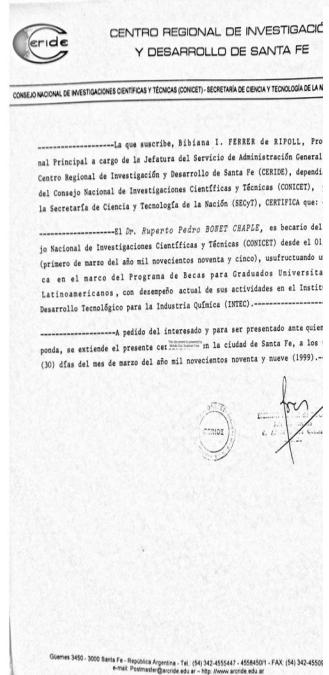
R. P. B.
RUBERTO PEDRO BONET CHAPLE

R. P. B.
PEDE R. E. HANCO
VERIFICADOR A/T DEL DOCUMENTO

J. F. P.
JUAN F. PÉREZITA
ADMINISTRADOR

J. F. P.
ROQUE M. R. ACOSTA COELLO
INVESTIGADOR
Registrado Universidad N° 49.362 2º folio

RESEARCH EXPERIENCES FELLOWSHIP PROGRAMS



AREA 15: INGENIERIA CIVIL			
 MINISTERIO DE CIENCIA Y TECNOLOGÍA		AREA 15: INGENIERIA CIVIL Y ARQUITECTURA	
APELLIDO 1	APELLIDO 2	NOMBRE	TITULO
Principal			
Presentación y artículo			
Características			
Relación de candidatos elegibles			
Oferta de los Centros			
Requisitos de los candidatos			
Funcionamiento			
Plazos			
Convocatoria			
Formularios			
Consultas más frecuentes			
Teléfonos y direcciones de consulta			
Clasificación temática ANEP			
1 PORTE	AGEL	FERNANDO	MODELOS DE PREVISION HIDROMeteorOLÓGICA EN BASE A LA MEDIDA DE LLUVIA EN EL PUNTO
2 CASTELLOTE	ARMEDO	MARTA MARIA	SIMULACION ACCELERADA DE LA LIXIVIACIÓN A LARGO PLAZO DE HORMIGONES DE ALTA CALIDAD PARA LAS PRESTACIONES DE LA CARACTERIZACION MECANOSTRUCTURAL Y MECANICO-QUÍMICA A LA APLICACION INGENIERIL
3 LOPEZ	GARELLO	CARLOS MARIA	MODELOS COMPUTACIONALES PARA ESTUDIOS HETEROGÉNEOS
4 PEREZ	APARICIO	JOSE LUIS	MECÁNICA COMPUTACIONAL DE MATERIALES AVANZADOS PARA APLICACIONES DE LA INGENIERIA
5 BONET	CHAPLE	RUPERTO PEDRO	MÉTODOS COMPUTACIONALES PARA LOS FENÓMENOS DE CONSOLIDACIÓN EN LA FLUIDODINAMICA DE ESTUARIOS, RÍOS Y COSTAS MECÁNICAS
6 SANJUAN	ALEXANDRE	ANA	SISTEMAS DE CONSOLIDACIÓN CON POLIMEROS PARA LA RECUPERACIÓN DE MATERIALES POROSOS UTILIZADOS EN OBRAS DE EDIFICACIONES
7 GONZALEZ	RODRIGUEZ	ERNESTO MAURICIO	MODELACION DE EVOLUCION MORFODINAMICA DE PLAYAS APLICABLE A PROCESOS DE MEDIO PLAZO

CONVOCATORIA CONCURSO

INICIO

© Ministerio de Ciencia y Tecnología

RESEARCH EXPERIENCES

RESEARCH STAGES



TU Delft University of Technology

Sr. Director: Amadeu Delshams
Departamento de Matemática Aplicada I
ETSEIB
Av. Diagonal, 647
08028 Barcelona
Spain

Your reference and date Our reference Office telephone Date
P.O. Box 5051 +31.15.2785530 13th March 2006

Subject Bonet Sub-division DIAM

Dear Director Delshams,

Prof. Rupert Pedro Bonet Chaple has visited our institute from February 15 - March 15, 2006. We have had a nice collaboration on the development of Finite Element Methods combined with stabilization techniques for the solution of advection-diffusion-reaction equations. During our weekly discussions, together with my colleague Dr. F.J. Vermolen, we observe a lot of progress on these methods. After some adaptations we think that the progress report, written by Prof. Bonet, can appear in our TU Delft Technical Report series.

In order to combine these techniques with our Level Set approach to solve moving boundary problems some more time and effort is needed. We hope that Prof. Bonet can stay another period in our institute in the near future to carry out this research.

Yours sincerely,

Dr.ir. C. Vulk

TECHNISCHE
UNIVERSITÄT
DRESDEN

Electrical Engineering, Mathematics and Computer Science

Fakultät für Mathematik und Naturwissenschaften
Der Direktor
Institut für Numerische Mathematik
Prof. Dr. Hans-Görg Roos

Technische Universität Dresden, 01062 Dresden
Sr. Director: Amadeu Delshams
Departamento de Matemática Aplicada I
ETSEIB
Av. Diagonal 647

08028 Barcelona
Spain

Dresden, 8. Februar 2006

Dear Director Delshams,

In our research group the Prof. Rupert Pedro Bonet Chaple have worked in February 1st - 9th period. This work with the PhD. Torsten Linß was related to stabilization of the advection-diffusion equations. On this period the professor Rupert teach a seminar entitled "Computational Methods to transport Equations in Coastal Process."

Your sincerely

Prof. Dr. Hans-Görg Roos
Direktor des Institutes für
Numerische Mathematik

Postadresse (Briefe): TU Dresden
D-01062 Dresden
Bauschule, Postfach 90050
Helmholtzstraße 10
01089 Dresden

Reisebüro Pakete u. d.:
TU Dresden
D-01062 Dresden

Besucheradresse:
Secretariat:
Zentraler Weg 12-14
Wittenau, D-01124

Zufahrt:
Hausnummern
Wittenau-C-Pflug
Internet: <http://www.tu-dresden.de>

BENEMÉRITA UNIVERSIDAD
AUTÓNOMA DE PUEBLA



FACULTAD DE CIENCIAS
FÍSICO MATEMÁTICAS

"Año de la Autonomía Universitaria"

DR. AMADEU DELSHAMS
DIRECTOR DPTO. MATEMATICA APPLICADA I
ETSEIB-UPC
PRESENTE:

Por este medio envío a usted un saludo cordial y al mismo tiempo extiendo una certificación de que el Profesor Dr. Rupert Pedro Bonet Chaple miembro de su Departamento Docente, ha visitado y trabajado en colaboración con los investigadores del Cuerpo Académico de Ecuaciones Diferenciales y Modelación Matemática de esta Facultad, en el periodo comprendido del 13 al 22 de Mayo del 2006.
En este periodo el Profesor Bonet participó en varias charlas y seminarios de discusión en los temas de colaboración conjunta con el grupo de investigación, también impartió un seminario sobre problemas objetos de investigación con los cuales él está vinculado. El profesor Bonet ha participado directamente en el desarrollo de un modelo matemático para la descripción de la dinámica del ciclo hidrológico en una cuenca. En su primera etapa se obtuvo la deducción y orientación de un modelo simplificado para el escurreimiento de la superficie libre con evaporación, transpiración y filtración. En este periodo el profesor trabajó en la orientación de un estudiante de la maestría en Matemáticas de esta facultad.

En espera de que en el futuro se den este tipo de colaboraciones en beneficio de nuestros estudiantes e investigadores, me es grato quedar de usted,

ATENTAMENTE

"PENSAR BIEN, PARA VIVIR MEJOR"
M. Puebla de Z., 22 de mayo de 2006.

DRA. ESPERANZA GUZMAN OVANDO
SECRETARIA DE INVESTIGACIONES
Y ESTUDIOS DE POSGRADO

"2006. AÑO DEL BICENTENARIO DEL NATALICIO DEL BENEMÉRITO DE LAS
AMÉRICAS DON BENITO JUÁREZ GARCÍA"

Av. San Claudio y Río Verde, Col. San Manuel, Ciudad Universitaria, Puebla, Pue., C.P. 72570
Tel.: (01 222) 229 55 00 Ext. 7550, 7552, Directo: (01 222) 229 56 37, Fax: (01 222) 229 56 36

When responding please quote our reference

Departament de Matemàtica
Aplicada I

UNIVERSITAT POLITÈCNICA DE CATALUNYA

Escola Tècnica Superior
d'Enginyeria Industrial de Barcelona
Diagonal, 647
E-08034 Barcelona
Tel. 93 401 65 48
Fax 93 401 17 13

Yo, Yuri Fedorov Kuzmin, el miembro del equipo de investigación del proyecto
MTM2009-12672,

Certifico, que recibí el investigador Rupert Pedro Bonet Chaple en el
departamento de Matemática Aplicada I de la UPC, durante el periodo
3 de febrero -- 19 de diciembre de 2014, y que hemos realizado juntos con el un
trabajo sobre temas del proyecto mencionado (estudio de singularidades
complejas de soluciones de algunos sistemas integrables y descripción de las
variedades invariantes complejas).

UPC
Departament
de Matemàtica Aplicada I

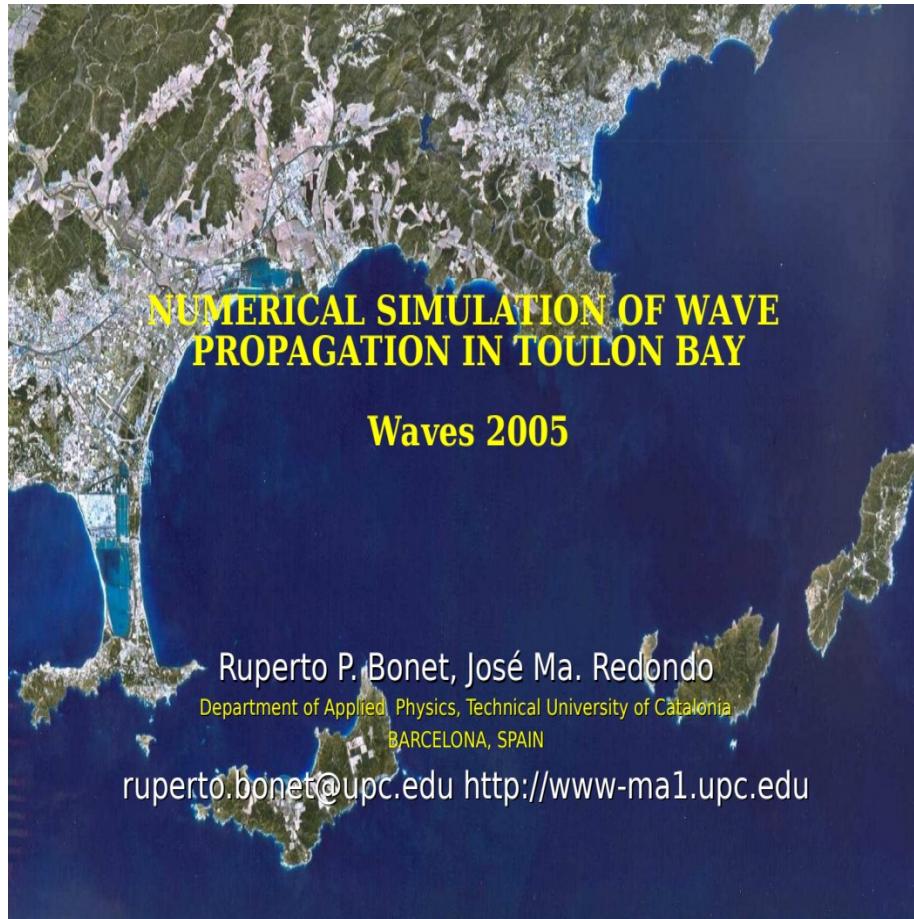
Barcelona, 30 de abril de 2015

Yuri Fedorov Kuzmin

23/01/2017

RESEARCH EXPERIENCES

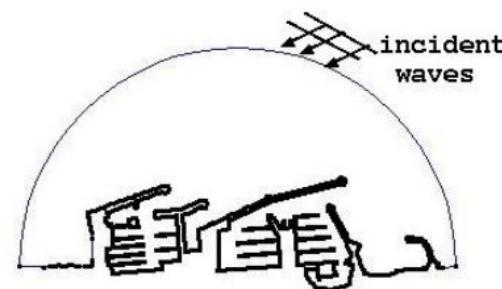
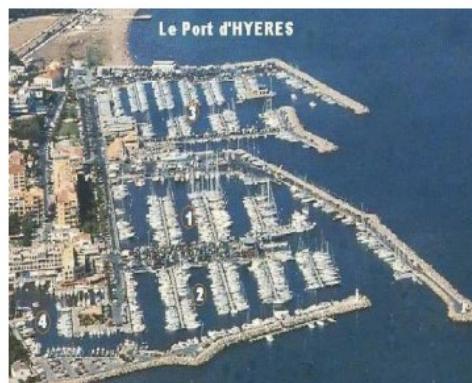
SIMULATION PROJECTS OF ENGINEERING PROBLEMS



RESEARCH EXPERIENCES

SIMULATION PROJECTS OF ENGINEERING PROBLEMS

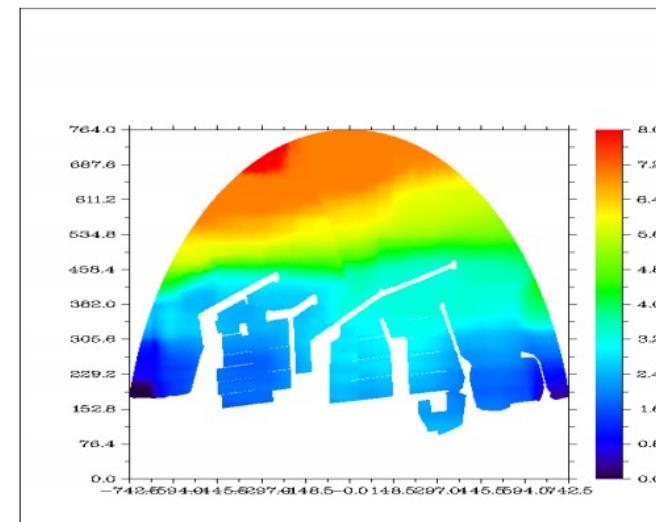
Waves in Hyeres Harbor



Period : 6s

Incidence angle : South-East

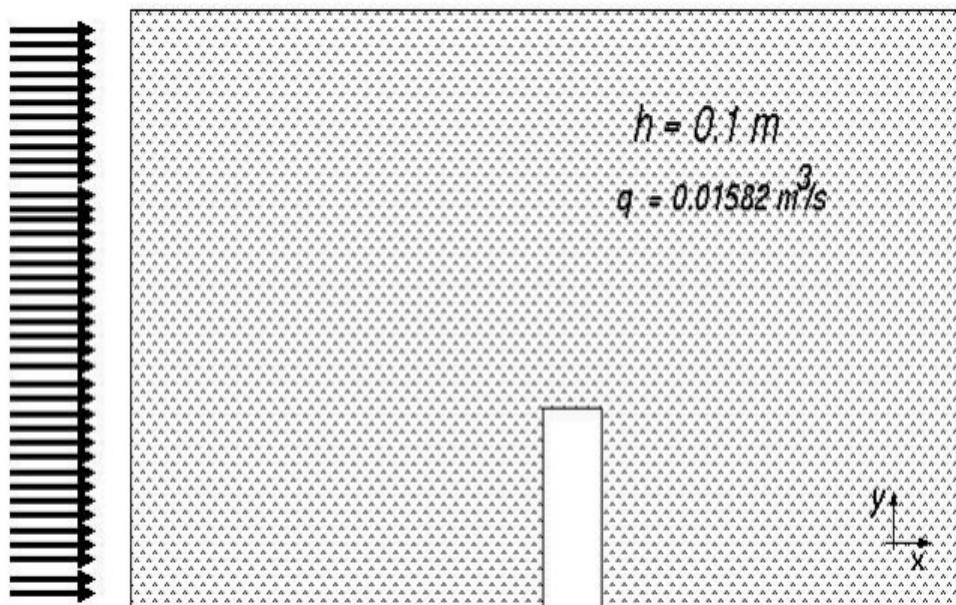
Fully Reflecting at coastline



Skills: Data entry, data preparation, C++, MATLAB, modelling, FEM

SPECIFIC PHYSICAL PROBLEMS

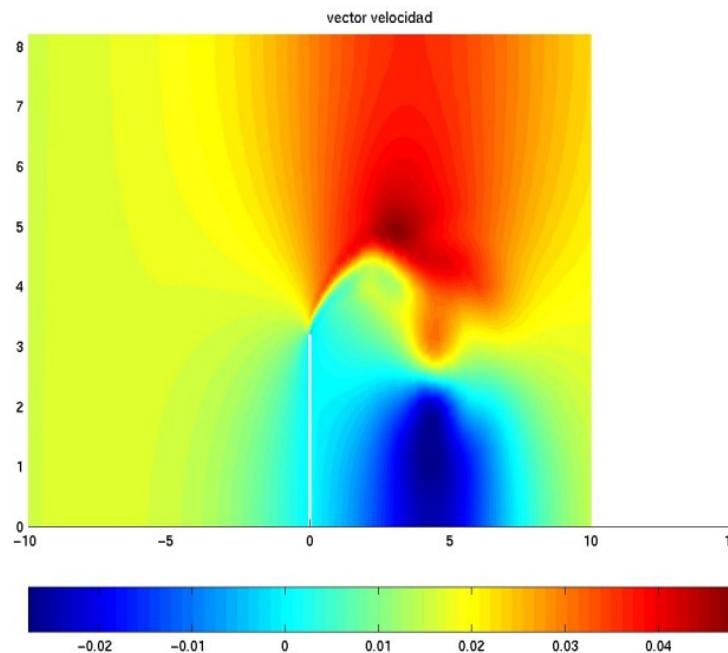
INVISCID FLOW IN OPEN CHANNELS WITH LATERAL CONTRACTION



RESEARCH EXPERIENCES

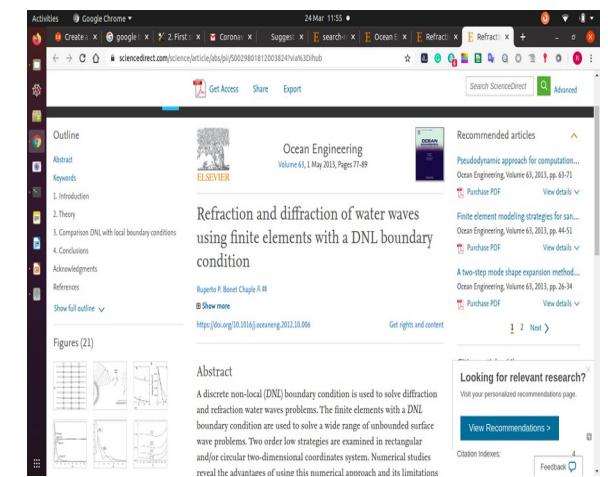
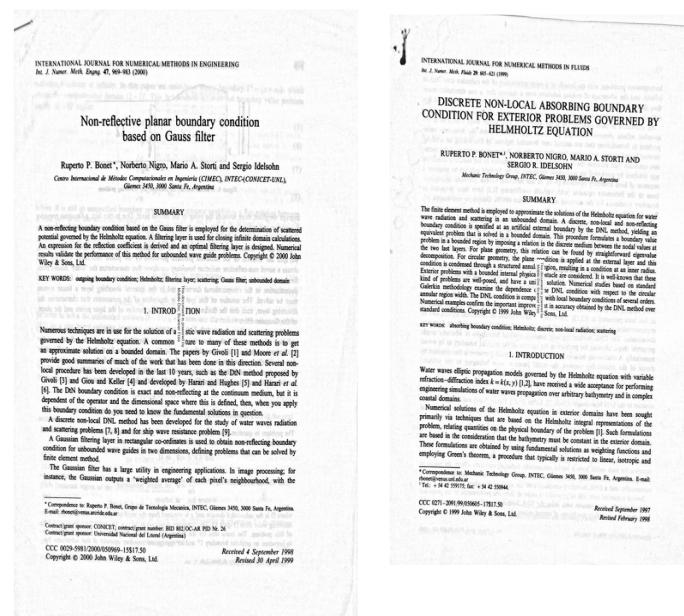
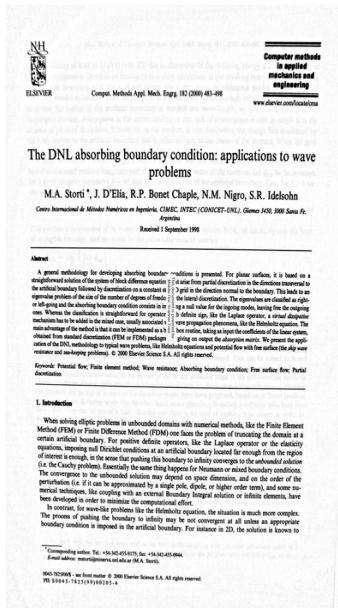
SIMULATION PROJECTS OF ENGINEERING PROBLEMS

VELOCITY FLOW MAP



Skills: Data entry, data preparation, C++, MPI, MATLAB, modelling, FEM

PUBLISHINGS



EXPERIENCES in STATISTICS



The screenshot shows a LibreOffice Writer document titled "SYLLABUS ESTADISTICA INFERENCIAL ABRIL-SEPT2016a.docx". The document contains the following text:

UNIVERSIDAD TÉCNICA DE BABAHoyo

FACULTAD DE ADMINISTRACION, FINANZAS E INFORMATICA
ESCUELA DE CONTADURIA, AUDITORIA Y FINANZAS
SYLLABUS DE LA ASIGNATURA
ESTADISTICA INFERENCIAL

I. INFORMACIÓN GENERAL

1 Carrera:	Ingeniería en Contabilidad y Auditoria				
1 Modalidad:	Presencial				
1 Nombre de la asignatura:	Estadística Inferencial				
1 Período académico:	Abril 2016 - Setiembre 2016				
1 Código:	ICA - 4402				
1 Eje de formación:	Básico.				
1 Año/Semestre	Séptimo Semestre				
1 Número de créditos:	5 CREDITOS				
1 Horas semanales:	Horas Presenciales	4	Horas Autónomas	4	Prácticas de aplicación y experimentación de los aprendizajes
1 Total Horas Año/Semestre	Horas Presenciales	64	Horas Autónomas	64	Prácticas de aplicación y experimentación de los aprendizajes
1 Prerrequisito:	Estadística descriptiva				

EXPERIENCES in STATISTICS



UNIVERSIDAD TÉCNICA DE BABAHYOY
FACULTAD DE ADMINISTRACIÓN FINANZAS E INFORMATICAS
ESCUELA DE CONTABILIDAD Y AUDITORIA



SYLLABUS ESTADÍSTICA DESCRIPTIVA

I. INFORMACIÓN GENERAL

1.1.	Carrera:	CONTADURIA, AUDITORIA Y FINANZAS			
1.2.	Modalidad:	PRESENCIAL			
1.3.	Nombre de la asignatura:	ESTADÍSTICA DESCRIPTIVA			
1.4.	Periodo académico:	SEPTIEMBRE 2015-DICIEMBRE 2015			
1.5.	Código:	X			
1.6.	Eje de formación:	PROFESIONAL			
1.7.	Semestre:	MATUTINO			
1.8.	Número de créditos:	X			
1.9.	Horas semanales:	Horas Presenciales	3	Horas Autónomas	Prácticas de aplicación y experimentación de los aprendizajes
1.10.	Total Horas Semestre: 160	Horas Presenciales	64	Horas Autónomas	Prácticas de aplicación y experimentación de los aprendizajes 32
1.11.	Prerrequisito:	ALGEBRA			
1.12.	Correquisito:	ESTADISTICA INFERNACIONAL			
1.13.	Duración periodo lectivo:	19 semanas	Fecha inicio Septiembre 28/09/2015	Fecha culminación 26/11/ 2015	E-mail:
1.14.	Profesores responsables:	Ec. Martha Acosta Roby, MCA Dra. Enrique Díaz Chong, Msc. Ing. Héctor Crespo Caicedo PhD. Lic.Ruperto Pedro Bonet Chaple		mgaacosta@utb.edu.ec hcrespoc@utb.edu.ec rbonet@utb.edu.ec	

EXPERIENCES in STATISTICS



30/08/2016

UNIVERSIDAD TÉCNICA DE BABAHoyo
VICERRECTORADO DE INVESTIGACIÓN Y POSTGRADO
INSTITUTO DE INVESTIGACIÓN Y DESARROLLO



Certifica a:

Ing. Ruperto Bonett

Por haber participado como Capacitador en el PROGRAMA DE FORTALECIMIENTO DE CAPACIDADES DOCENTES EN INVESTIGACIÓN, en los Módulos de Concepción y Escritura de Artículos Científicos, MATLAB, e-learning/m-learning/b-learning, Matemática y Tecnología, Metodología de la Investigación, Diseño de Experimentos y Estadísticas en la Investigación durante los días los días 21 y 28 de octubre, 4 y 11 de noviembre y 2 de diciembre de 2015, en la Universidad Técnica de Babahoyo.



Dr. Rafael Falconí Montalván, MSc.
RECTOR

Lic. Adelita Pinto Yerovi, MSc.
VICERRECTORA DE INVESTIGACIÓN
Y POSTGRADO

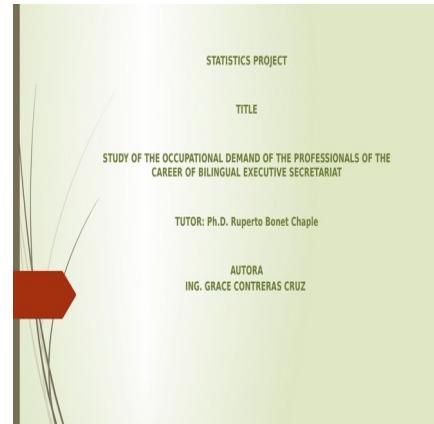
Babahoyo, 2 de diciembre de 2015

EXPERIENCES in STATISTICS (ADVISED PROJECTS)



ADDITION OF
EXOGENOUS
ENZYMES IN THE
FEEDING OF
BROILER CHICKENS

CARMEN VÁSCONEZ
MONTÚFAR
2016



UNIVERSIDAD TÉCNICA DE BABAHOYO

APPLIED STATISTICS COURSE

TITLE: Evaluation of the use of the Augmented Reality technology tool in the students of First and Second Semester of the Social Communication Career in the period April - September 2016.

NOMBRE: Silvia Paulina Maldonado Mangui

PROFESOR: Phd. Ruperto Bonet

Babahoyo - Ecuador
2016

TOPIC
MARKETING OF ORNAMENTAL PLANTS ON THE INTERNET

Docente
Ph.D. Ruperto Bonet Chaple

Autora
Msc. Ángela Jordán Y.

UNIVERSIDAD TÉCNICA DE BABAHOYO

ALTERNATIVES FOR THE CONTROL OF
MELOIDOGYNE spp AND RADOPHOLUS SIMILIS IN
TROPICAL FLOWERS

Ing.Agr. Emma Lombeida García
MBA



TECHNICAL UNIVERSITY OF BABAHOYO
Research and Development Institute
Applied Statistics Course

TITLE

The teaching performance and its influence on the educational quality of the Faculty of Legal, Social and Educational Sciences, of the Technical University of Babahoyo, Babahoyo canton, Los Ríos province, period September 2015 / February 2016.



Autora: Lcda. Glenda Intriago A.

Tutor: Dr. Ruperto Bonet C.

AI and BIG DATA



COURSERA CERTIFICATIONS

CLOUD COMPUTING

The image displays three side-by-side screenshots of cloud computing platforms:

- Google Cloud Platform API & Services dashboard:** Shows traffic and error metrics for a Python Data Dashboard project. It includes a sidebar for APIs & Services and a Cloud Shell terminal window.
- IBM Watson Studio projects list:** Displays a list of projects such as Natural Language, Tensorflow, Capstone Project Notebook, Apache SystemML, Scale a Keras model with IBM Watson, Run a Notebook using Keras and DL4J, Anomaly Detection, and WatsonDataFormV0.0.
- Microsoft Azure Solutions page:** Features a section titled "Azure solutions" with a sub-section "Find the solution to meet the needs of your application or business". It also includes a "Explore" section with links to Industries, Solution architectures, and Migration Center.

JAVA-HADOOP

The screenshot shows the Hadoop Web UI interface. At the top, there's a navigation bar with 'Activities' and 'Firefox Web Browser'. Below it, a tab bar shows 'sednabn/Python-JAVA-X' and 'RUNNING Applications - Mozilla Firefox'. The main content area is titled 'RUNNING Applications'. On the left, there's a sidebar with icons for Cluster, Applications, Scheduler, and Tools. The 'Cluster' section shows metrics like 'Apps Submitted' (0), 'Apps Pending' (0), 'Apps Running' (0), and 'Memory Used' (0 B). The 'Applications' section lists stages: NEW (1), SUBMITTED (0), ACCEPTED (0), RUNNING (0), FINISHED (0), FAILED (0), and KILLED (0). The 'Scheduler' section shows 'Scheduler Metrics' with columns for Active Nodes, Decommissioning Nodes, Decommissioned Nodes, Lost Nodes, and Unhe. The 'Tools' section shows a table for 'Scheduler Metrics' with columns for ID, User, Name, Application Type, Queue, Priority, StartTime, FinishTime, State, FinalStatus, Running Containers, Allocated CPU, Allocated Vcores, and Allocated Memory. A note at the bottom says 'Showing 0 to 0 of 0 entries'.

The Project Gutenberg EBook of Ulysses, by James Joyce

This eBook is for the use of anyone anywhere at no cost and with almost no restrictions whatsoever. You may copy it, give it away or re-use it under the terms of the Project Gutenberg License included with this eBook or online at www.gutenberg.org

Title: Ulysses

Author: James Joyce

Release Date: August 1, 2008 [EBook #4300]

Last Updated: October 30, 2018

Language: English

Character set encoding: UTF-8

*** START OF THIS PROJECT GUTENBERG EBOOK ULYSSES ***

Produced by Col Choat, and David Widger.

cover

Ulysses

by James Joyce

Contents

- I -

[1]
[2]
[3]

- II -

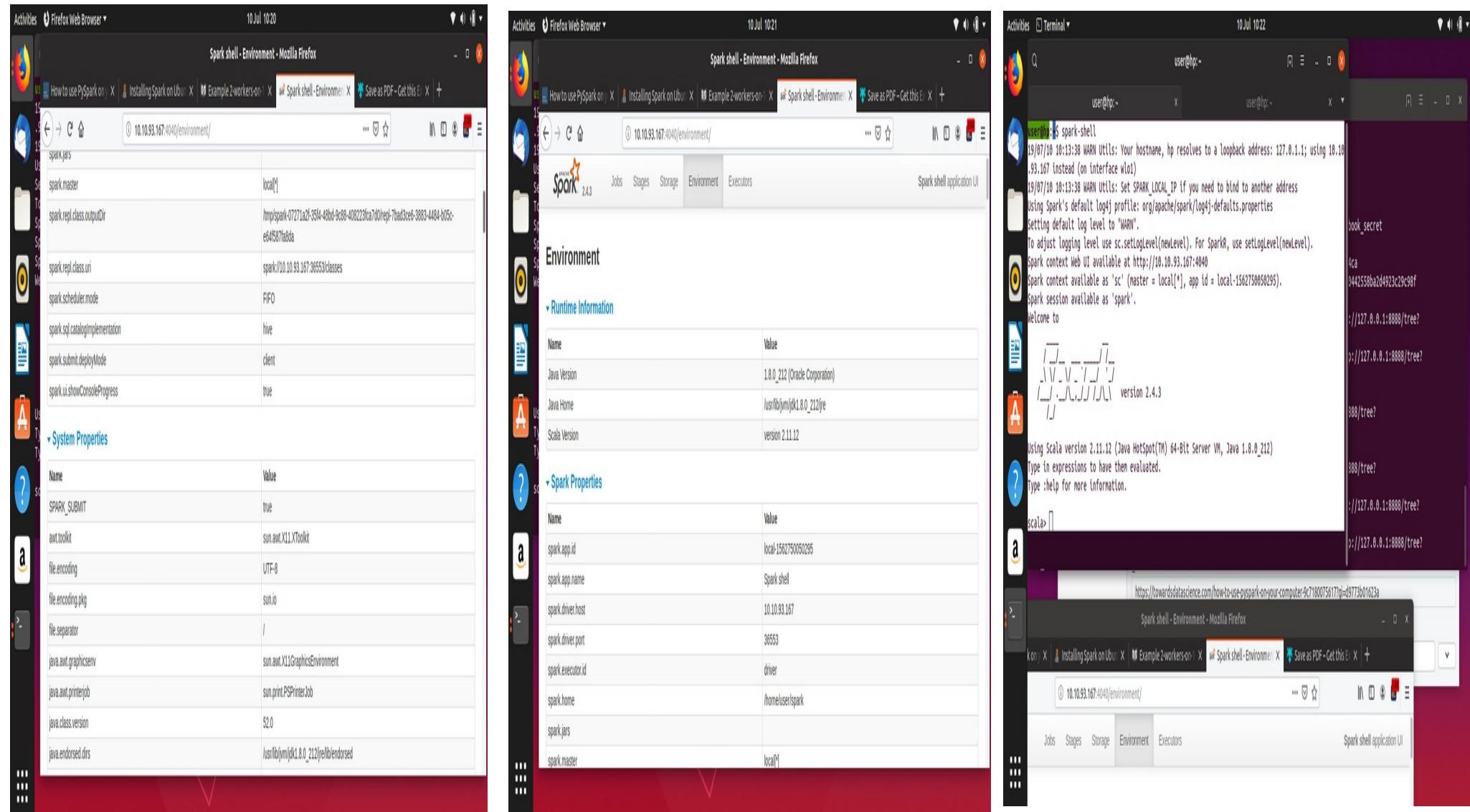
[4]
[5]
[6]
[7]
[8]
[9]
[10]
[11]
[12]
[13]
[14]
[15]

- III -

[16]

'(Lo)cura" 1
'1498 1
'1498," 1
'35" 1
'40, " 1
'A 2
'AS-IS". 1
'A_ 1
'Absoluti 1
'Aesopi" 1
'Alack! 1
'Alack!" 1
'Alla 1
'Allegorical 1
'Alpha 1
'Alpha," 1
'Alpine-glow" 1
'An 2
'And 3
'Antoni 1
'At 1
'BOILING" 2
'B_ 1
'Batesian" 1
'Bathers 1
'Beta 2
'Beta" 1
'Big 1
'Bononiae 1
'But 1
'By 1
'Cave-men" 2
'Clean 1
'Come 1
'Cromagnards" 1
'Crookes 2
'DARWIN'S 2
'Dagoes," 1
'Daily 1
'Day" 2
'De 1
'Death 1
'Defects". 1
'Defects," 2
'Description 1
'Disposizione 2
'Doctrinal 1
'E 1
'Egli 1
'El 2
'Elements". 1
'Every 2
'Fifty-foot 1
'First 2
'Florentie 1
'For 2
'Gamma 1
'Gamma" 1
'Georgics" 1
'Go 1
'Guido 2
"Here 1
"Historia 1
"How 1

SPARK



EXPERIENCES AS A BIG DATA ANALYTICS TUTOR

01/04/2020

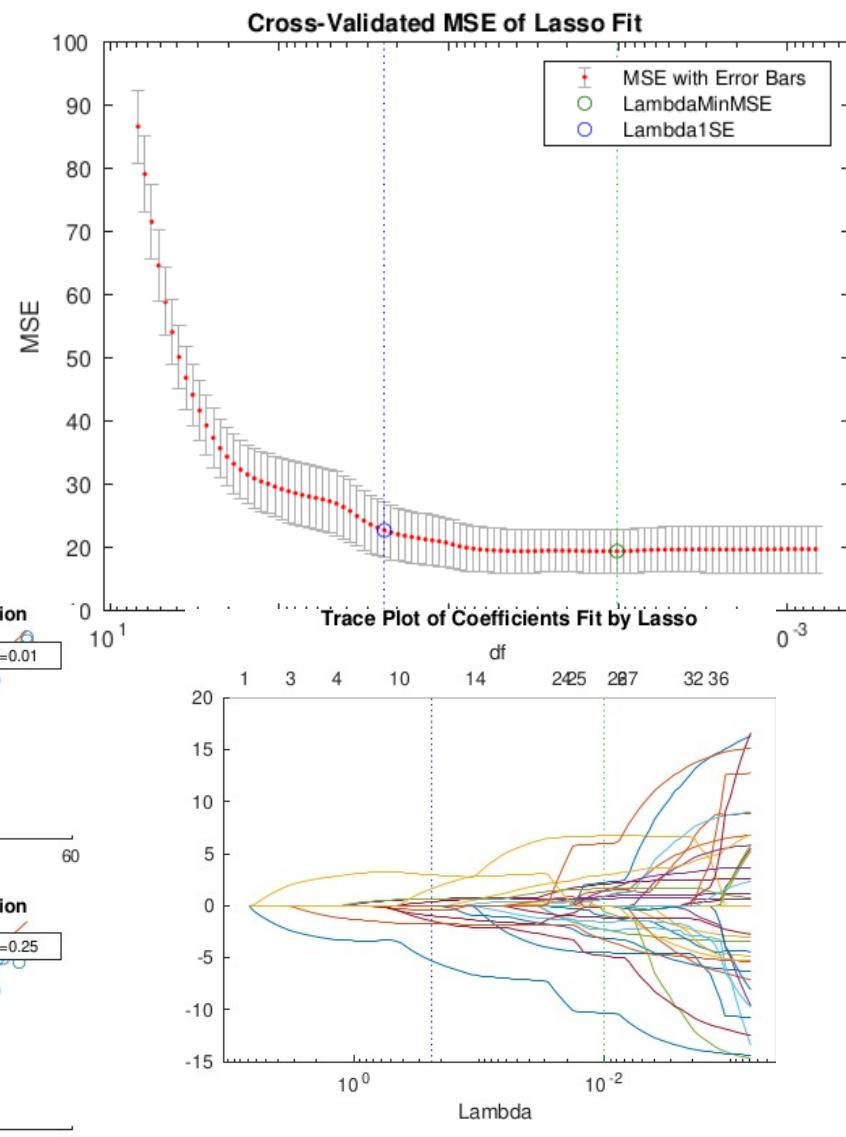
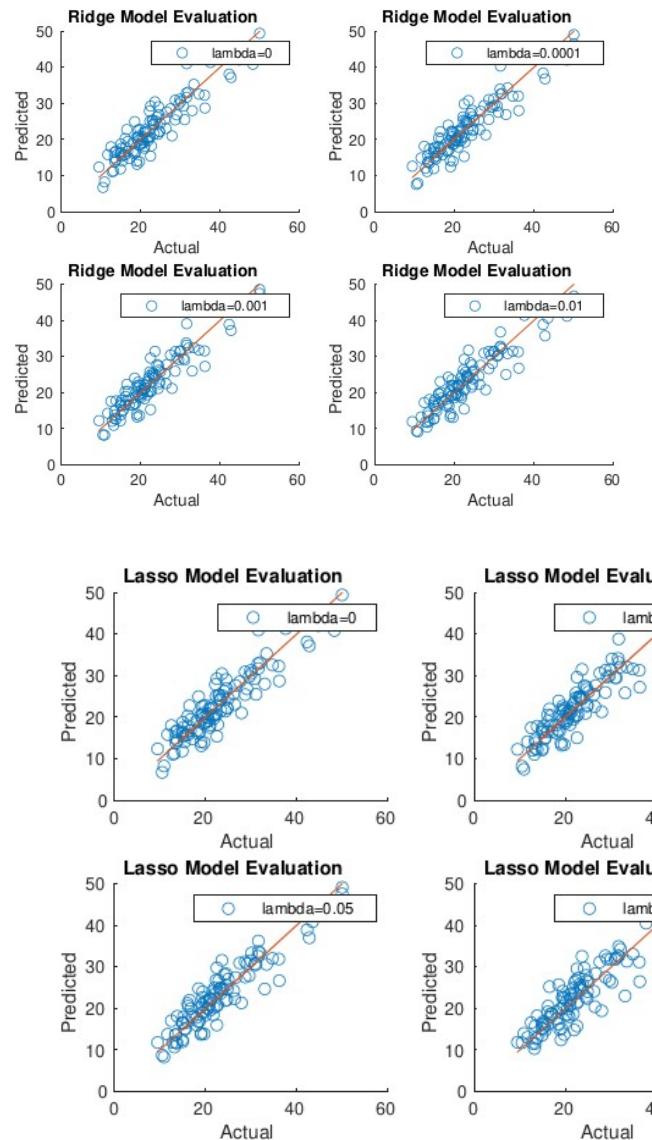
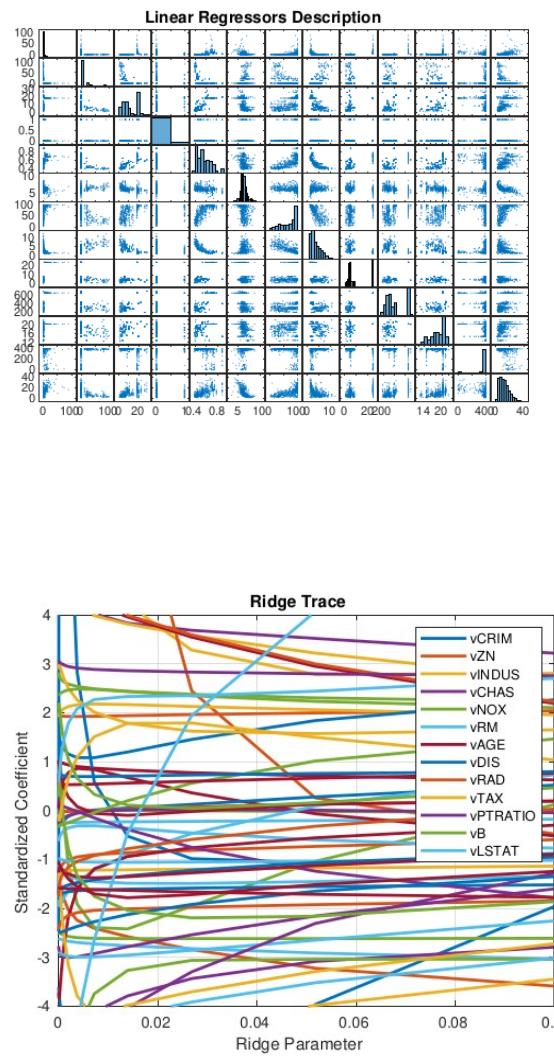
Copy_of_pset1_q2.ipynb - Colaboratory

```
1  clear all;
2  clf;
3  %Boston housing data from: http://lib.stat.cmu.edu/datasets/boston
4  %originally from a paper by Harrison and Rubinfield (1978)
5
6  % Variables in order:
7  % CRIM per capita crime rate by town
8  % ZN proportion of residential land zoned for lots over 25,000 sq.ft.
9  % INDUS proportion of non-retail business acres per town
10 % CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
11 % NOX nitrogen oxides concentration (parts per 10 million)
12 % RM average number of rooms per dwelling
13 % AGE proportion of owner-occupied units built prior to 1940
14 % DIS weighted distances to five Boston employment centres
15 % RAD index of accessibility to radial highways
16 % TAX full-value property-tax rate per $10,000
17 % PTRATIO pupil-teacher ratio by town
18 % B 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
19 % LSTAT % lower status of the population
20 % MEDV Median value of owner-occupied homes in $1000's
21
22 %read in the data, and define the 14 originally observed variables:
23
24 data = dlmread('housing.csv', ' ', 0,0);
25
26
27 vCRIM = data(:,1);
28 vZN = data(:,2);
29 vINDUS = data(:,3);
30 vCHAS = data(:,4);
31 vNOX = data(:,5);
32 vRM = data(:,6);
33 vAGE = data(:,7);
34 vDIS = data(:,8);
35 vRAD = data(:,9);
36 vTAX = data(:,10);
37 vPTRATIO = data(:,11);
38 vB = data(:,12);
39 vLSTAT = data(:,13);
40 vMEDV = data(:,14);
41
42 %the outcome variable is Median value of owner-occupied homes in $1000's:
43
44 y=vMEDV; %outcome variable
45 n=length(y); %number of observations
46
47 %define all the regressors:
48
49 X0=[vCRIM,vCRIM.^2,vCRIM.^3,... %cubic in crim (per capita crime rate)
50 vZN>0,... %dummy for zn > 0
51 vZN,vZN.^2,vZN.^3,... %cubic in zn (proportion residential land zoned for lots over 25,000 ft^2)
52 vINDUS,vINDUS.^2,vINDUS.^3,... %cubic in indus (proportion of non-retail business acres)
53 vCHAS,vCHAS,... %Charles River dummy
54 vNOX,vNOX.^2,vNOX.^3,... %cubic in nox (nitrogen oxide concentration)
55 vRM,vRM.^2,vRM.^3,... %cubic in rm (average rooms per dwelling)
56 vAGE,vAGE.^2,vAGE.^3,... %cubic in age (proportion of owner-occupied units build pre 1940)
57 vDIS,vDIS.^2,vDIS.^3,... %cubic in dis (weighted mean of distances to 5 employment centers)
58 vRAD,vRAD.^2,vRAD.^3,... %cubic in rad (index of highway accessibility)
59 vTAX,vTAX.^2,vTAX.^3,... %cubic in tax (property tax rate per $10,000)
60 vPTRATIO,vPTRATIO.^2,vPTRATIO.^3,... %cubic in ptratio (pupil-teacher ratio by town)
61 vB,vB.^2,vB.^3,... %cubic in black (1000*(proportion black - .63)^2)
62 vLSTAT,vLSTAT.^2,vLSTAT.^3,... %cubic in lstat (percent lower SES)
63 ];
64
65 p=size(X0,2); %total number of regressors (excluding the intercept here)
66
```

https://colab.research.google.com/drive/1OP2xh6XRY-sj6C2Cq7YdjUWNelsU_YG3#printMode=true

EXPERIENCES AS A BIG DATA ANALYTICS TUTOR

Boston Housing Dataset



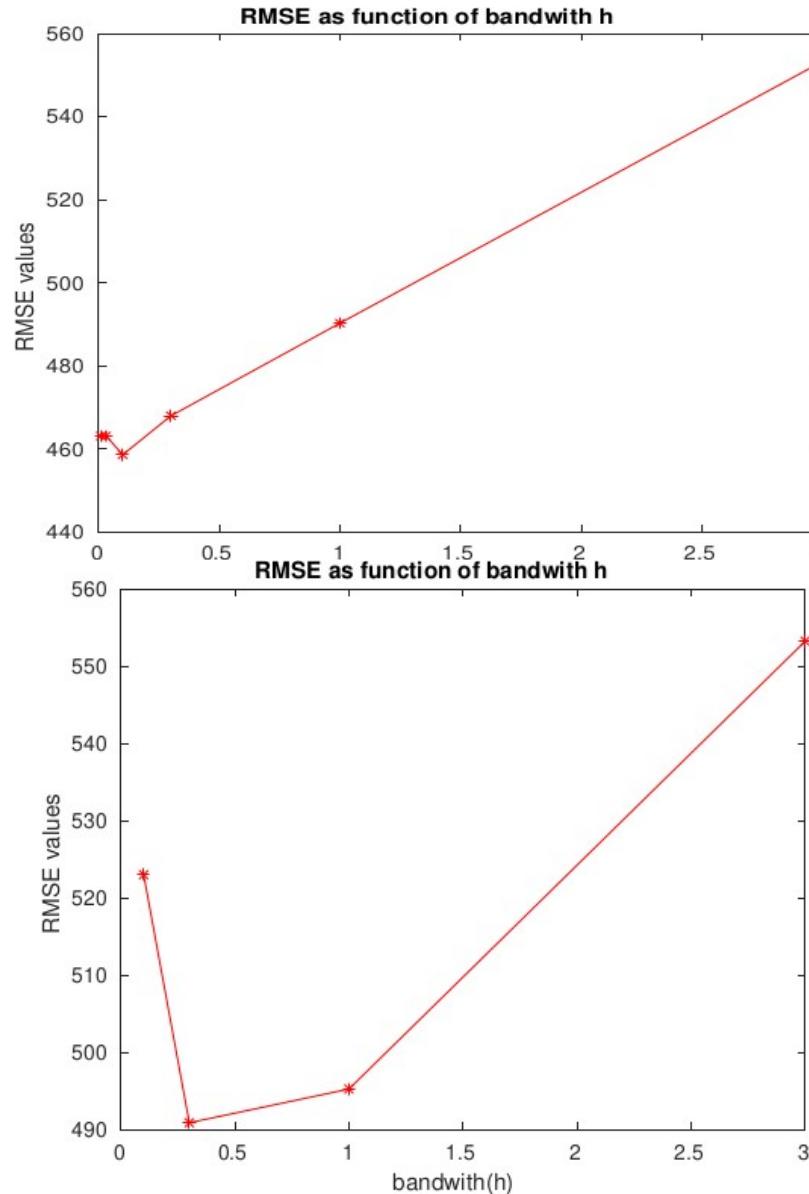
EXPERIENCES AS A BIG DATA ANALYTICS TUTOR

Smoking Effect on Birth Outcomes

```

1 clear all;
2 %DATA FROM J. Abrevaya, "Estimating the Effect of Smoking on Birth Outcomes Using a
3 %Matched Panel Data Approach," Journal of Applied Econometrics, Vol. 21,
4 %No. 4, 2006, pp. 489-518
5 %
6 %DATA DOWNLOADED FROM: http://qed.econ.queensu.ca/jae/2006-v21.4/abrevaya/
7 %
8 %THE DATA HAVE ALREADY BEEN RANDOMLY PERMUTED
9 %
10 %
11 %READ DATA FROM FILE
12 %
13 %
14 data = dmatread('birpanel.csv', ' ', 0);
15 names={'momid'; 'dx'; 'states'; 'dmage'; 'dmeduc'; 'mrbirth'; 'gestar'; 'dbirth'; 'cigar'; 'smoker'; 'male'; 'year'; 'married'; 'hsgrad'; 'somecoll'; 'collgrad'; 'agesq'; 'black'; 'idecode2'; 'idecode3'; 'novail'; 'pretil2';
16 ind=fin(data, 10,-99);
17 data=data(ind,:);
18 %
19 %DROP OBSERVATIONS FOR WHICH VARIABLE 'cigar' IS UNKNOWN
20 %
21 %
22 ind=fin(data, 10,-99);
23 data=data(ind,:);
24 %
25 %DEFINE RELEVANT VARIABLES FOR OUR EXERCISE:
26 %
27 dmage=data(:,4); %age of mother (in years)
28 dmeduc=data(:,5); %education of mother (in years)
29 gestar=data(:,6); %length of gestation (in weeks)
30 agesq=(data(:,7)).^2; %square of age of mother
31 agesq=(data(:,7)).^3; %cube of age of mother
32 agesq=(data(:,7)).^4; %fourth power of age of mother
33 cigar=(data(:,10)); %number of cigarettes smoked per day (99=unknown)
34 %
35 %DEFINE OUTCOME VARIABLE = birthweight (in grams)
36 %
37 %
38 y=dbirth;
39 %
40 %DEFINING COVARIATES:
41 %
42 % age of mother, education of mother,
43 % square of age of mother, number of cigarettes smoked per day
44 %
45 %
46 x=[mage dmeduc gestar cigar];
47 %
48 %Standardize all the covariates to have std = 1
49 %
50 %
51 %
52 %
53 for k=1:4
54 x(:,k) = x(:,k) / std(x(:,k));
55 end
56 %
57 %
58 %Decompose dataset into training set and validation set:
59 %
60 %
61 n=length(y); %total sample size
62 %
63 Y=y(1:(n-100));
64 X=x(1:(n-100));
65 nT=length(Y); %sample size training set
66 %
67 Y_val=y((n-100):(n));
68 X_val=x((n-100):(n));
69 nV=length(Y_val); %sample size validation set
70 %
71 %
72 %QUESTION 1%
73 %
74 h_vector=[0.0, 0.03, 0.1, 0.3, 1];
75 [ypred_vector,rmse_vector]=kernel_regression_vector(h_vector,nT,nV,xT,yT,xV,yV);
76 %
77 figure(1)
78 plot(h_vector,rmse_vector,'*')
79 xlabel('bandwidth(h)')
80 ylabel('RMSE values')
81 title('RMSE as function of bandwidth h')
82 disp(h_vector)
83 disp(rmse_vector)
84 %
85 %
86 %QUESTION 2%
87 %
88 RMSE_OLS=rmse_vector(1);
89 beta=xt'\xt; %xt'xt;
90 xT_av=mean(xT,:);
91 yT_av=mean(yT,:);
92 beta_OLS=xt\av_beta;
93 y_OLS=beta_0 + xT\beta;
94 y_prediction=beta_0 + xV\beta;
95 rmse_OLS=sqrt(mean_square_error,y);
96 disp(RMSE_OLS)
97 %
98 %
99 %QUESTION 3%
100 nT=10000;
101 xt=X(1:nT,:);
102 yt=Y(1:nT);
103 [ypred_vector,rmse_vector]=kernel_regression_vector(h_vector,nT,nV,xT,yT,xV,yV);
104 %
105 figure(2)
106 plot(h_vector,rmse_vector,'*')
107 xlabel('bandwidth(h)')
108 ylabel('RMSE values')
109 title('RMSE as function of bandwidth h')

```



MACHINE LEARNING

SiMLeng-Statsmodels package

```

#!/usr/bin/env python3
from simulation_statsmodels import Simleng_strategies as Simleng
Simleng(0,-1).strategies()

#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Fri Sep 28 14:46:37 2018

@author: sedna
"""

import pandas as pd
import numpy as np
from collections import OrderedDict
import matplotlib.pyplot as plt
from matplotlib import colors as mcolors
from biokit.viz import corrplot
import seaborn as sns
from scipy import stats
import statsmodels.api as sm
import statsmodels.discrete.discrete_model as smd
from statsmodels.multivariate.pca import PCA as smPCA
from sklearn.decomposition import PCA as skPCA
from sklearn import metrics
from statsmodels.stats.outliers_influence import variance_inflation_factor
from metrics_classifier_statsmodels import MetricBinaryClassifier as ac
from tools import Tools
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis as QDA
colors = [ic for ic in mcolors.BASE_COLORS.values() if ic !=(1,1,1)]
class Data_base:

    def data_generation():
        """LOADING THE PIMA DATA SET"""
        data_train={}
        data_test={}
        pima_tr = pd.read_csv('pima.tr.csv', index_col=0)
        pima_te = pd.read_csv ('pima.te.csv', index_col=0)
        # Training aata
        df=pima_tr
        df.dropna(how="all", lace=True) # drops the empty line at file-end
        columns=df.columns
        columns_train=columns[7]
        columns_test=columns[-1]
        X_train=df[columns_train]
        Y_train=df[columns_test]
        # Testing data
        de=pima_te
        X_test=de[columns_train]
        Y_test=de[columns_test]
        data_train['train']=[columns_train,X_train,Y_train,df]
        data_test['test']=[columns_test,X_test,Y_test,de]
        return data_train,data_test

    def data_generation_binary_classification(**kwargs):
        """Generation of dummy variables to Binary Classification Task"""

        _X_train,_Y_train,_=Tools.data_extract_dict_to_list('train',**kwargs)
        _X_test,_Y_test,_=Tools.data_extract_dict_to_list('test',**kwargs)

        data_dummy_train={}
        data_dummy_test={}
        # Dummy variables to categorical variable
        V_train=Tools.data_dummy_binary_classification(_Y_train,"No",0,1)
        V_test=Tools.data_dummy_binary_classification(_Y_test,"No",0,1)

```



```

#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Mon Sep 17 10:05:57 2018

@author: sedna
"""

#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Sat Sep 29 12:03:20 2018

@author: sedna
"""

class Tools:

    def data_unpack_kwargs(**kwargs):
        # data is a list of keys and values
        # kind='list' to list
        #kind='dict' to dict
        keys=kwargs.keys()
        data={}
        for name in keys:
            data[name]=kwargs[name]
        return data

    def data_list_unpack_dict(dict):
        data=[]
        for name in dict.keys():
            data_list=dict[name]
            if len(data_list)==1:
                data_list=list(dict[name])[0]
            data.extend(data_list)
        return data

    def data_extract_dict_to_list(kind,**kwargs):
        dataDict={}
        for k,name in enumerate(kwargs.keys()):
            if name==kind:
                data_dict[kind]=data[kind]
                data_list=Tools.data_list_unpack_dict(data_dict)
                return data_list

    def data_list_to_matrix(x,shape):
        import numpy as np
        Lx=len(x)
        Ls=len(shape)
        if Ls==2: [Nx,Ny]=shape[:]
        if Ls==3: [Nx,Ny,Nz]=shape[:]
        if Ls==2:
            matrix=np.empty((Nx,Ny))
            for ii in range(Lx):
                matrix[:,ii]=x[ii]
            assert Nx,Ny==matrix.shape
        return matrix

    def data_list_to_transform_to_integer(x):
        """get a new type for the members of the list """
        xx=[]
        for i in x:
            xx.append(int(i))
        return xx

    def add_data_list_of_several_lists(nn,x):
        LEN=0
        for ii in range(nn):
            LEN+=len(x[ii])
        return LEN,print(LEN)
        def add_data(x):
            import numpy as np
            n,m=x.shape
            x=np.array(x)
            sum_by_rows=[x[ii,:].sum() for ii in range(n)]
            return np.asarray(sum_by_rows).astype(float)

```

SiMLeng PROJECT



PIMA INDIAS DATASET

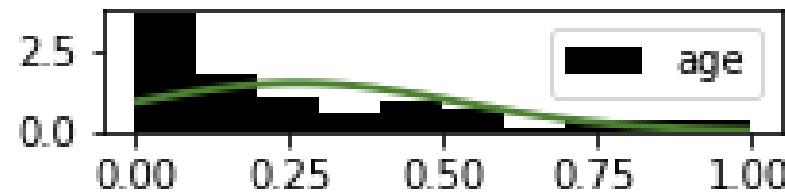
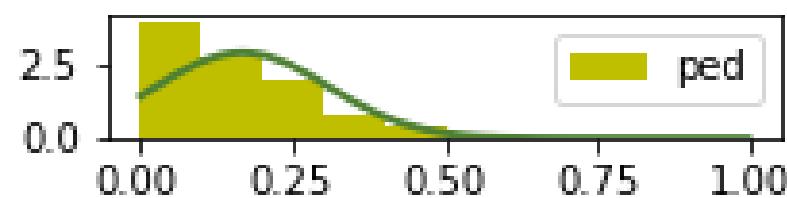
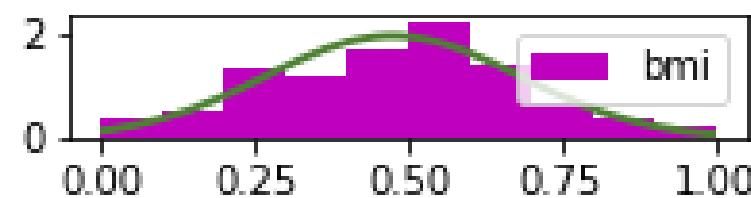
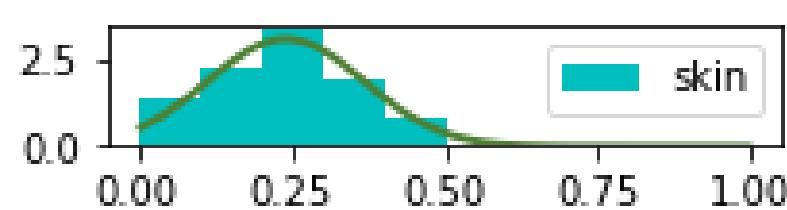
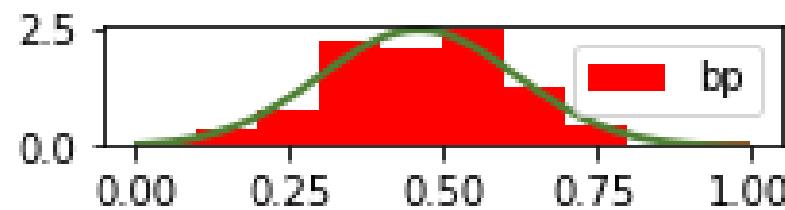
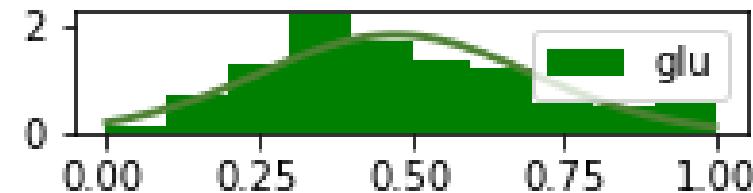
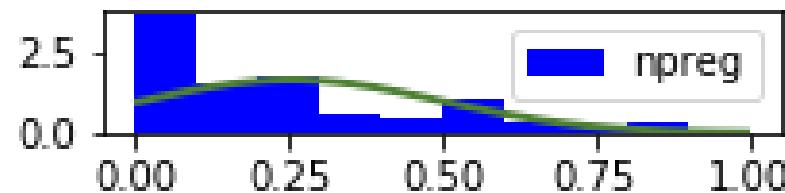
CLASSIFICATION TRAINING

RUPERTO P. BONET

rpbonetch@gmail.com

Predictors Histogram-Normal pdf

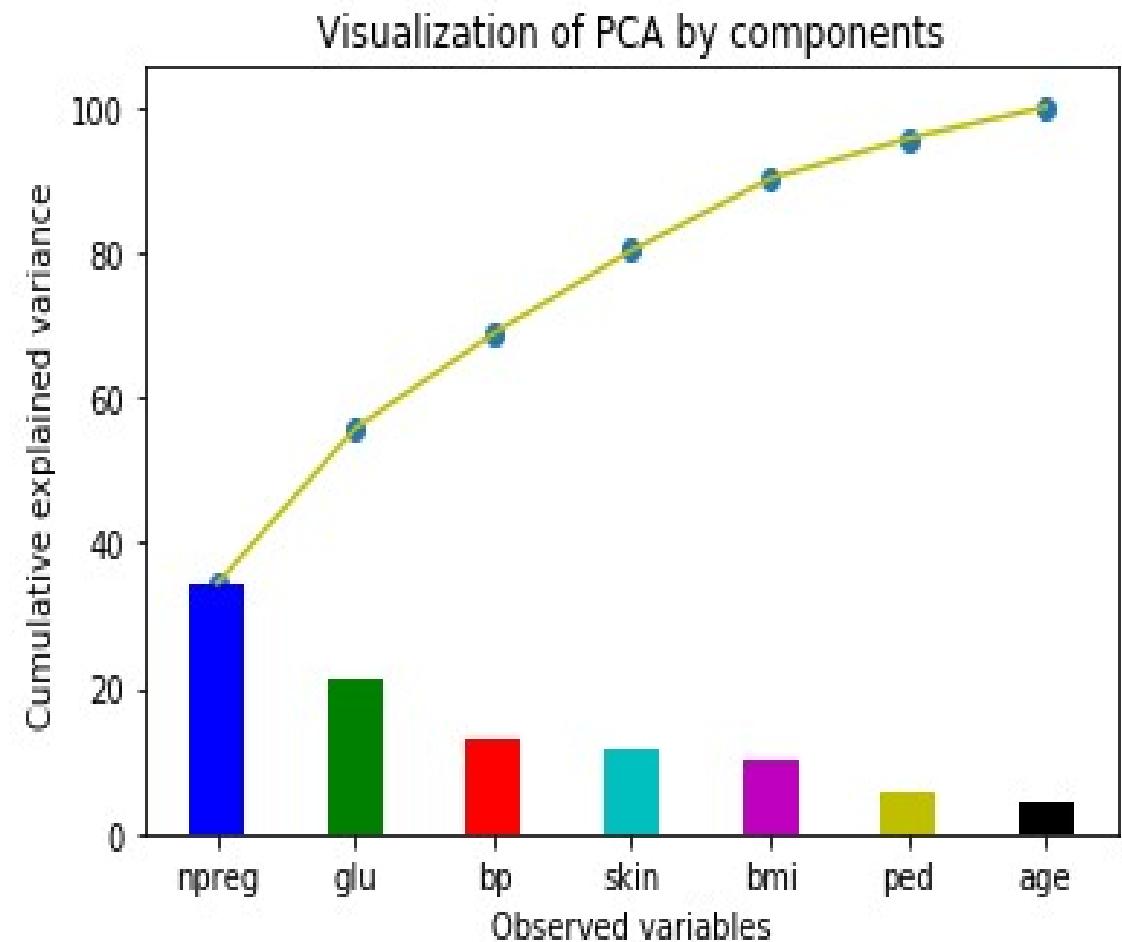
Pima Datasets [type~npreg+glu+bp+skin+bmi+ped+age]



FEATURES-SELECTION

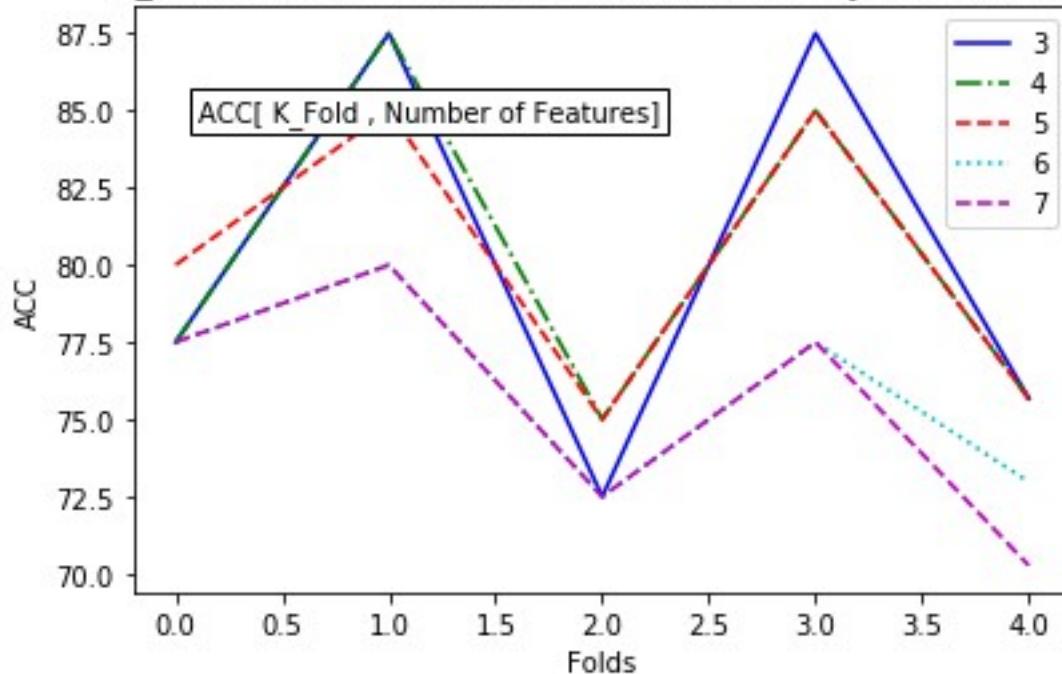
PCA (ncomp=7)

Eigenvalues in descending order			
	ods_var	Eig.Abs	cum_var_exp
0	npreg	481.852	34.418
1	glu	299.289	55.796
2	bp	182.374	68.823
3	skin	159.994	80.251
4	bmi	138.027	90.110
5	ped	11.964	95.679
6	age	60.500	100.000

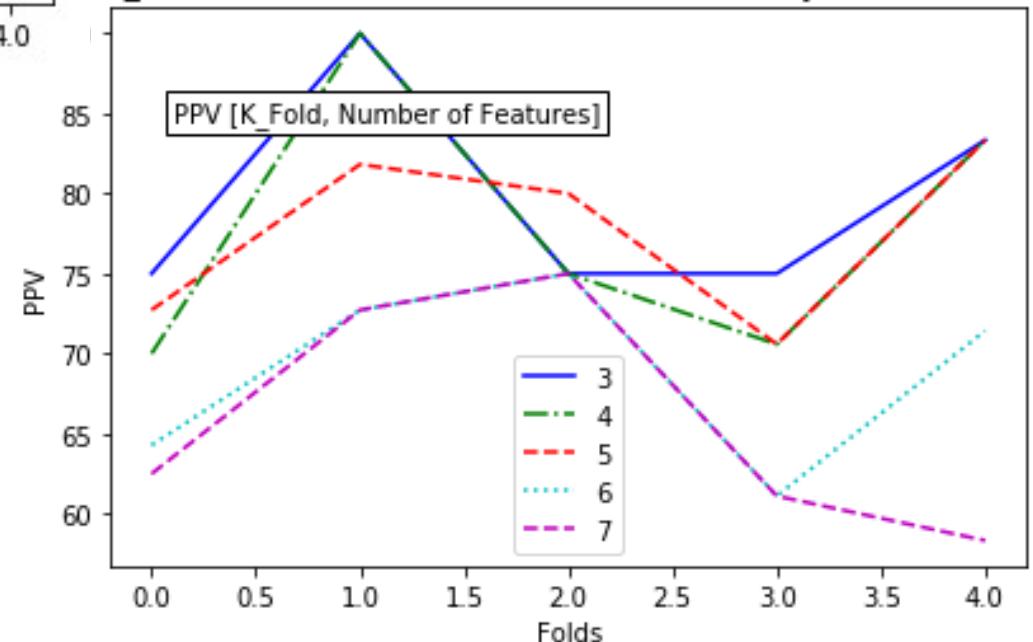


CROSS VALIDATION K_FOLD=5

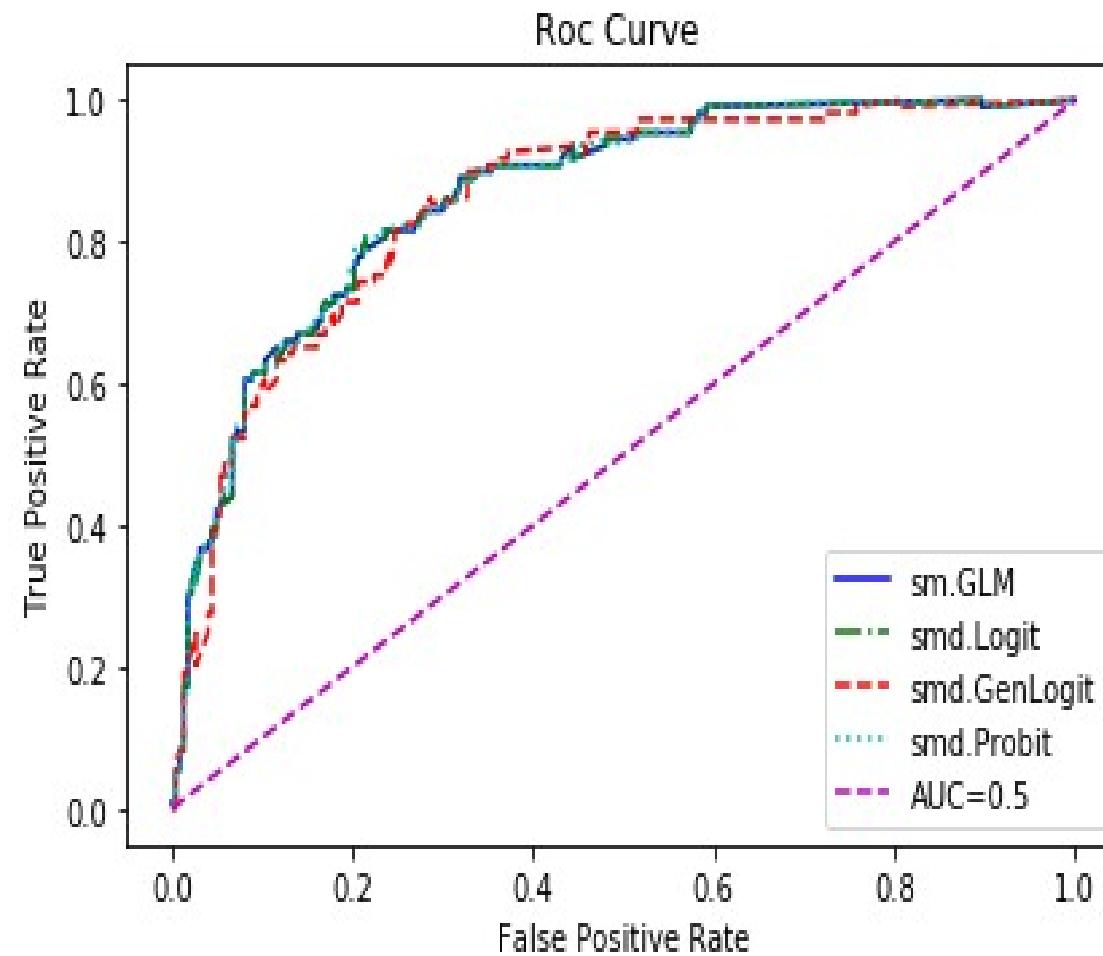
K_Fold vs Features : Cross-Validation to Binary Classification



K_Fold vs Features : Cross-Validation to Binary Classification



BINARY CLASSIFICATION RESULTS



PRELIMINARY CONCLUSIONS

- Glu and ped are most relevant predictors.
- Genlogistic function with c lower than 0.5 has strong influence on the detection of TRUE NEGATIVE, and by consequence, it improves the prediction performance.
- The best prediction with 3 predictors is 81.93%, which has been obtained using GenLogistic function (with c=0.2) on the features subset: npreg, bmi and ped.
- The best prediction with 4 predictors is 80.12%.
- This prediction on the features subset: npreg, glu, bmi and ped has been obtained applying two different strategies:
 - QDA to the original variables, and
 - Logit to the transformed variables by PCA with four components
- This result is consistent for the study and this features subset represents 63.22% from training data information.

SiMLeng PROJECT

Genlogistic vs Logistic

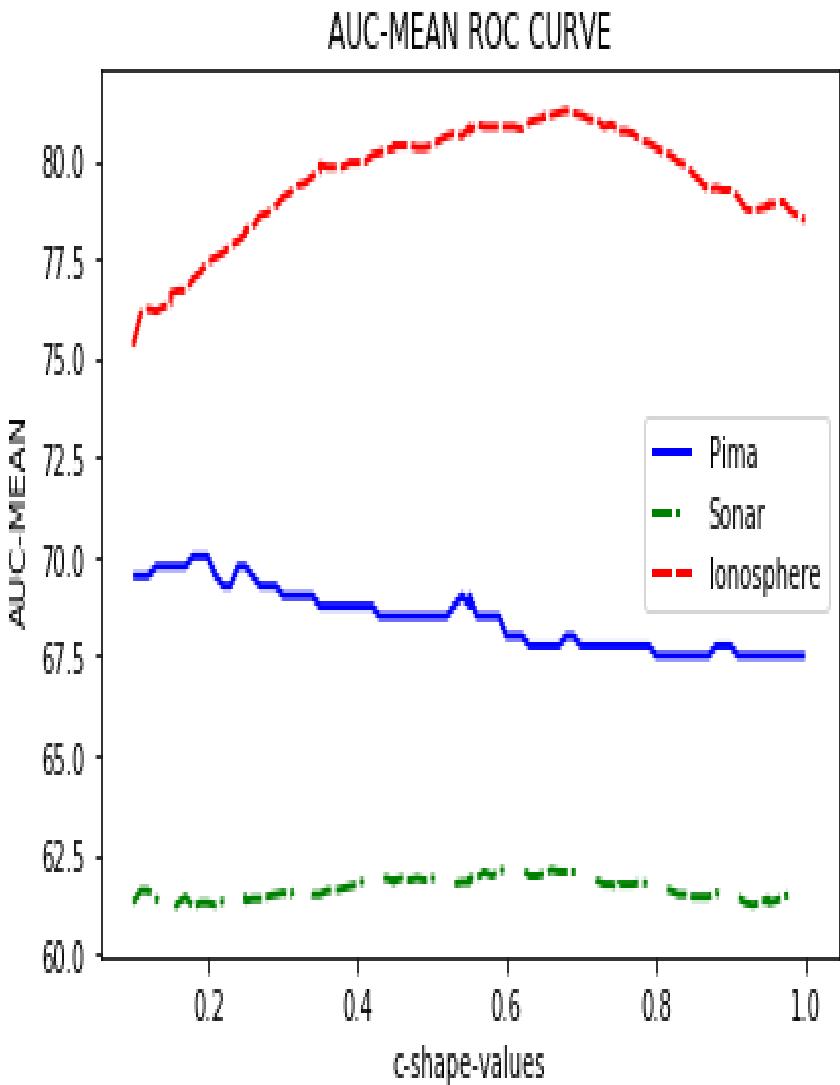
**NOTE ON
BINARY CLASSIFICATION**

RUPERTO P. BONET

rpbonetch@gmail.com

COMPARISON BETWEEN DATASETS

AUC MEAN(%)



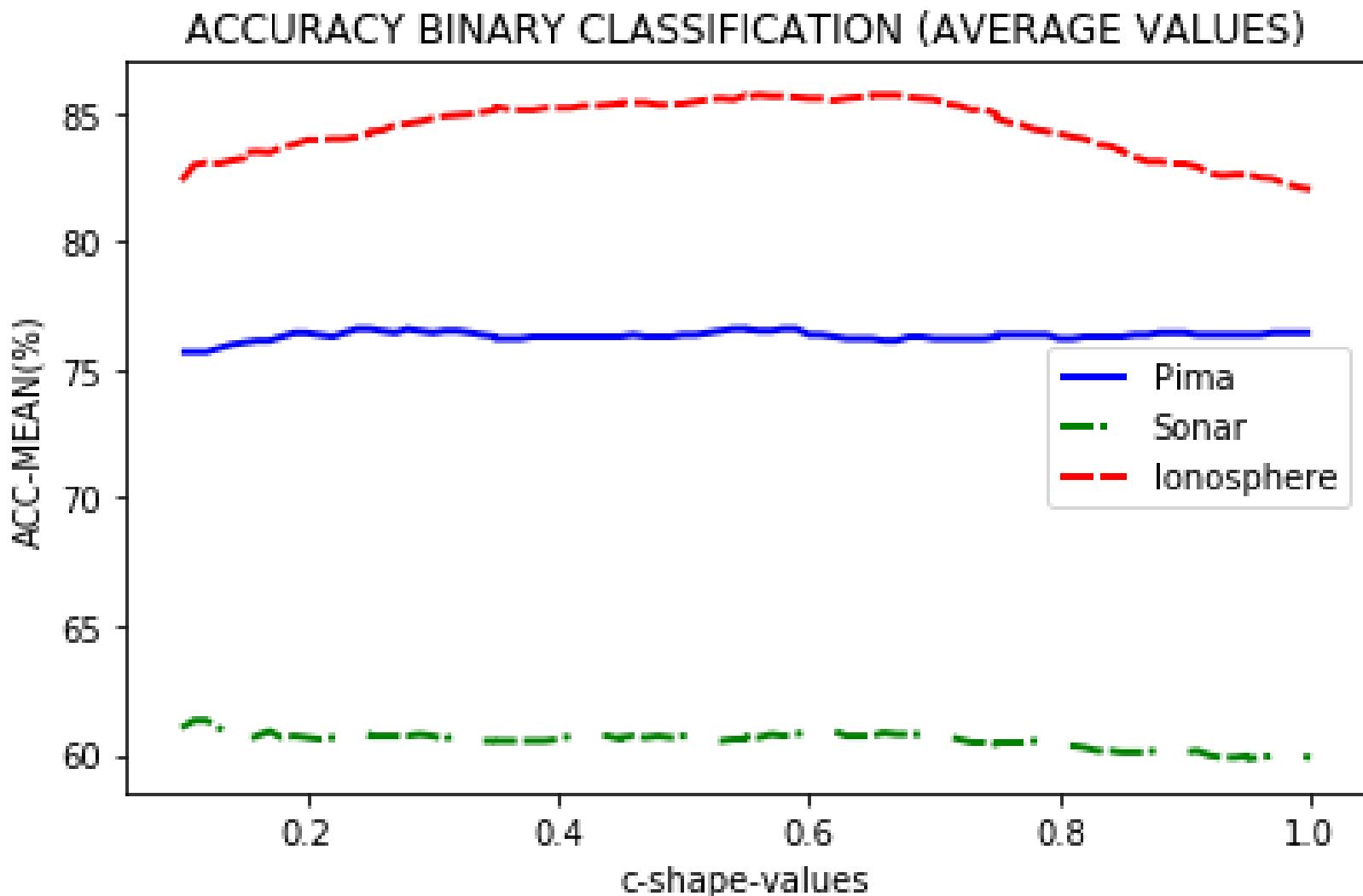
```

1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3
4 Created on Sun Nov 4 22:17:43 2018
5 These classes are added to statsmodels package
6
7 @author: sedna
8
9 import numpy as np
10 from statsmodels import base
11 from statsmodels.stats.discrete.discrete_model import BinaryModel
12
13 class GenLogit(BinaryModel):
14     """
15     Binomial choice genlogit model
16     """
17     %(%(params)s
18     %(extra_params)s)
19
20     Attributes
21     -----
22     endog : array
23         A reference to the endogenous response variable
24     exog : array
25         A reference to the exogenous design.
26     *** : {params : base.model_params, doc,
27            'extra_params' : base.missing_param_doc}
28
29     def get_c_value(self):
30         self.c = fcat(input("Shape genlogistic: "))
31         self.c=.35
32         self.c=0.2
33         return self.c
34
35     def cdf(self, X):
36         """
37         The genlogistic cumulative distribution function
38         """
39         Parameters
40         -----
41         X : array-like
42             X is the linear predictor of the logit model. See notes.
43         c : float
44             c is a shape parameter c>0 |c|=1 becomes to logistic case.
45         Returns
46         -----
47         1/(1+exp(-|x|)*c)
48
49     Notes
50     -----
51     In the genlogit model,
52
53     .. math:: \text{Lambda}(\left| x \right|) \prime | \beta_0, c | = \ln(\text{Prob}) / \left| x \right|
54
55     .. math:: \text{pdf}(\text{self}, X)
56
57     c = self.get_c_value()
58     X = np.asarray(X)
59     return 1/(1+np.exp(X))**c
60
61     def pdf(self, X):
62         """
63         The logistic probability density function
64
65         Parameters
66         -----
67         X : array-like
68             X is the linear predictor of the logit model. See notes.
69         c : float
70             c is a shape parameter
71         Returns
72         -----
73         pdf : ndarray
74             The value of the Logit probability mass function, PMF, for each
75             observation in X given the parameters X1, ..., Xn.

```

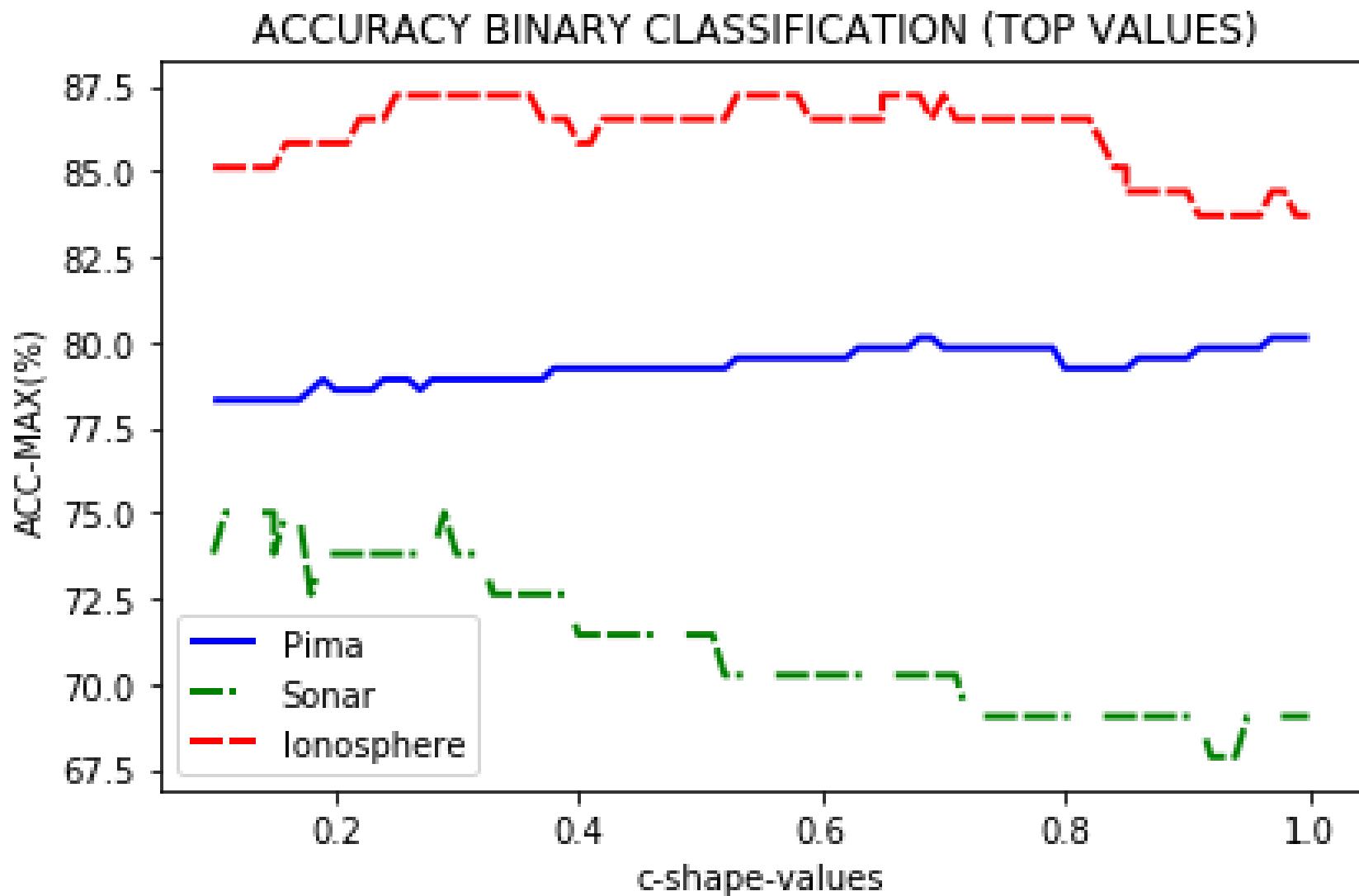
COMPARISON BETWEEN DATASETS

ACC MEAN(%)



COMPARISON BETWEEN DATASETS

ACC TOP VALUES(%)



PRELIMINARY CONCLUSIONS

Results seem to indicate that Genlogistic function could reduce the false positive/ negative rate more than Logistic predictions. The following table shows the above observations:

Datasets	Pima Indias			Sonar			Ionosphere		
TOP VALUES(%)	AUC	ACC	FPFN	AUC	ACC	FPFN	AUC	ACC	FPFN
Logistic (c=1)	73.0	80.12	19.88	70.0	69.05	30.95	81.0	83.69	16.31
Genlogistic($0 < c < 1$)	77.0	80.12	19.88	75.0	75.00	25.00	83.0	87.23	12.77

SiMLeng PROJECT

Lumped Sample Collocation Prediction Method

RUPERTO P. BONET

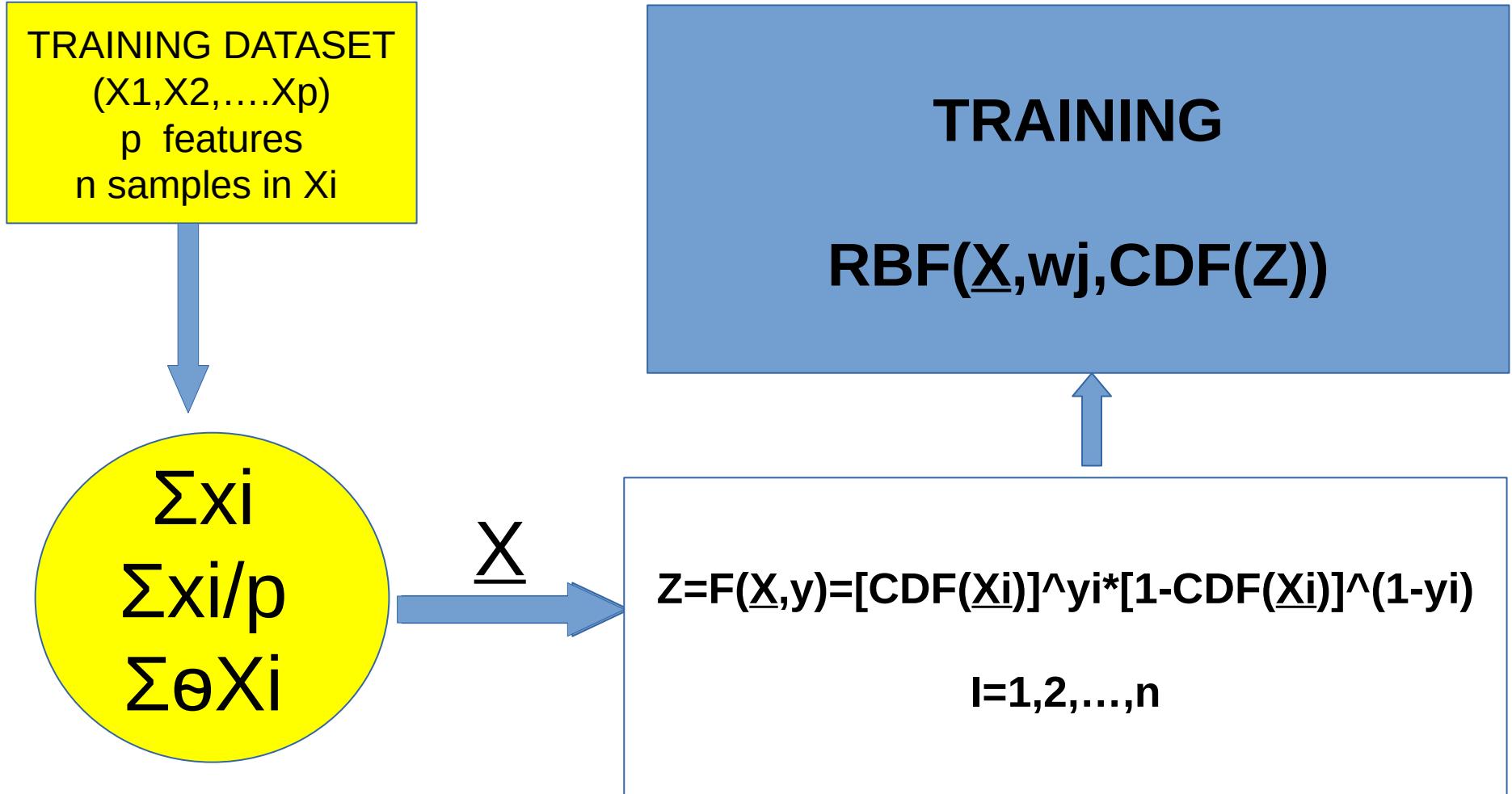
rpbонетч@gmail.com

INTRODUCTION

A numerical method based on PCA features selection and lumping transformations of samples in combination with One-Dimensional Radial Basis Functions Collocation method is implemented.

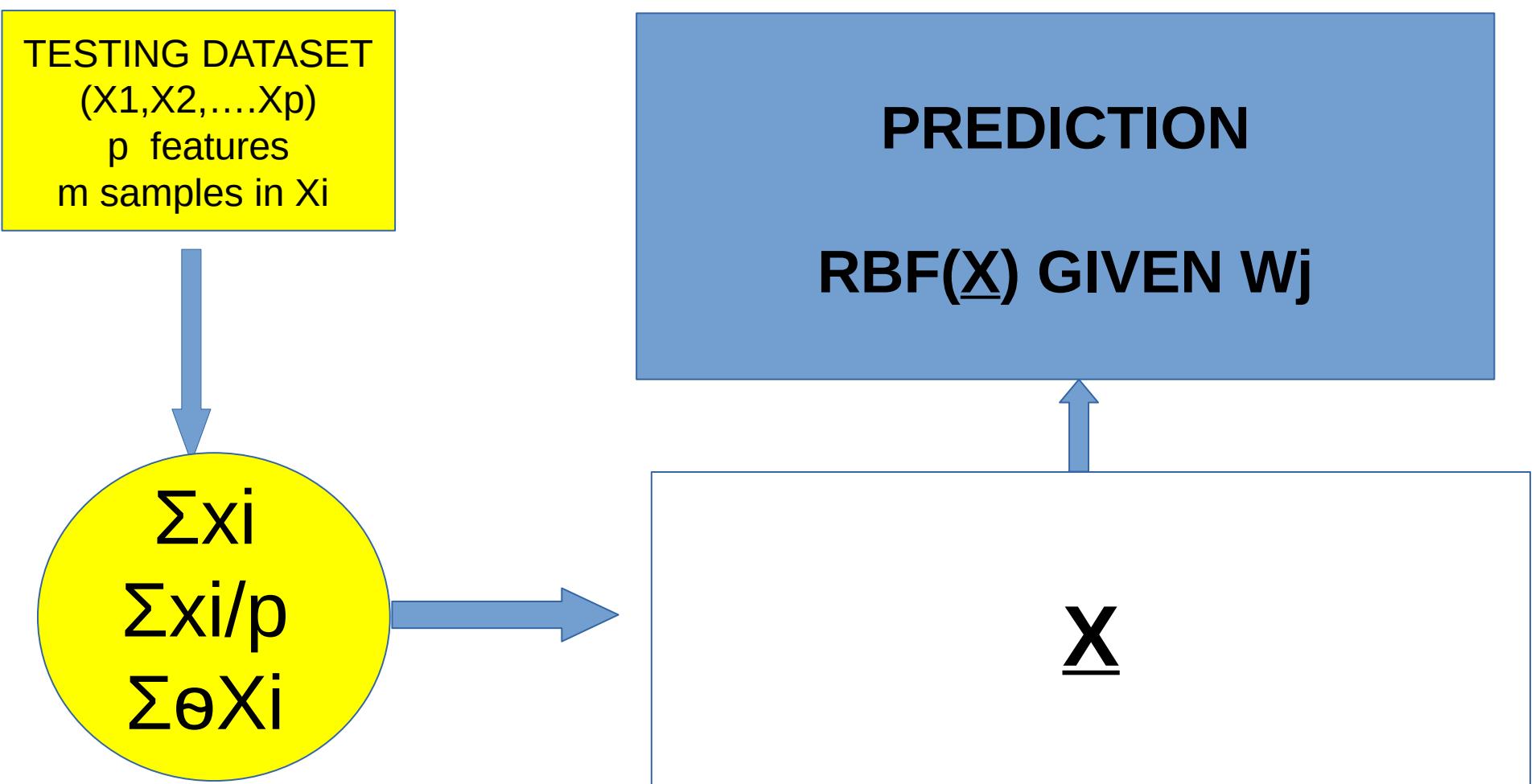
This experiment has been designed using Pima Indians, Sonar and Ionosphere datasets, which are related to making prediction and binary classification.

FEATURES-SAMPLE LUMPING COLLOCATION PREDICTION METHOD



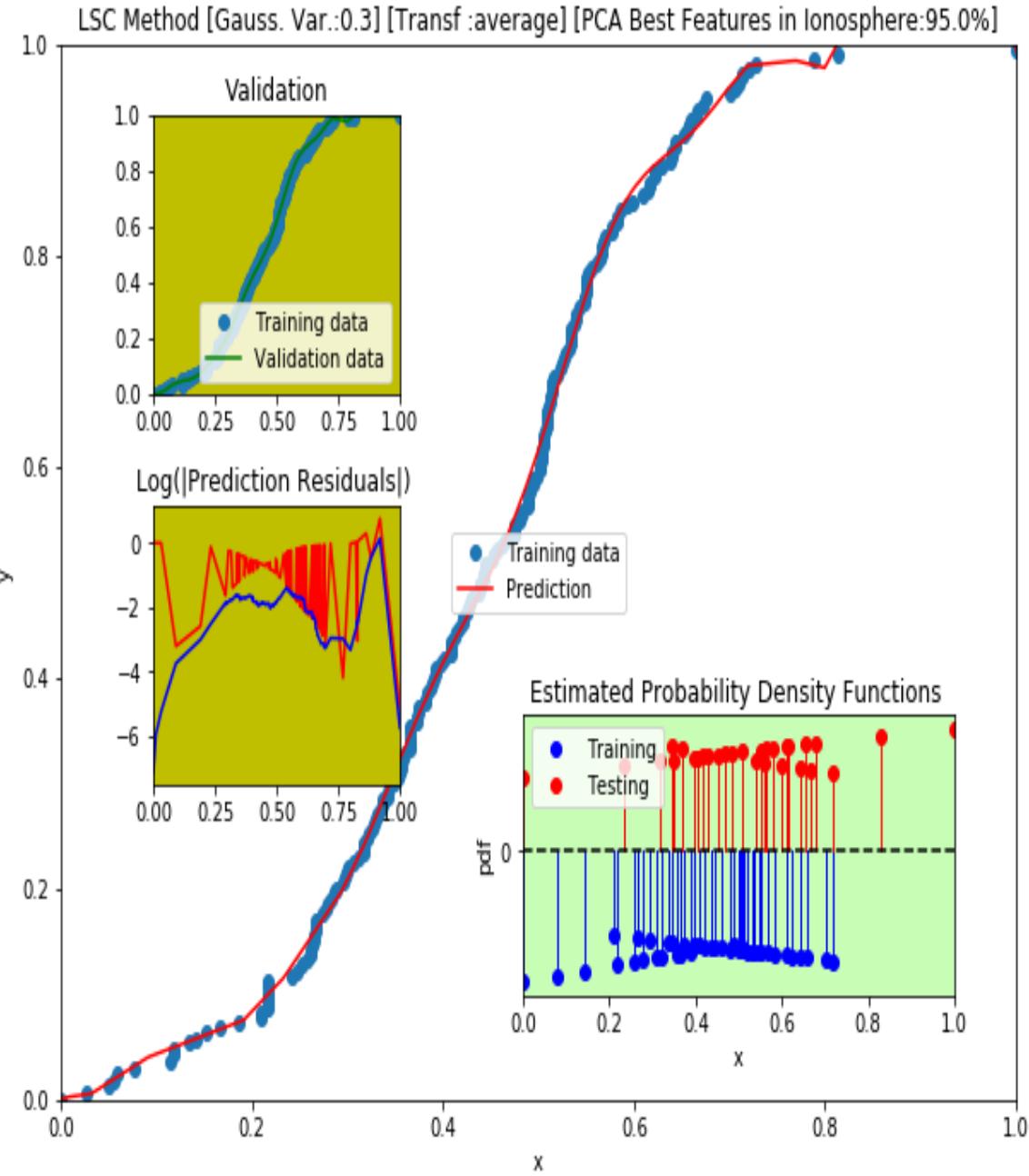
TRANSFORMATIONS

FEATURES-SAMPLE LUMPING COLLOCATION PREDICTION METHOD



TRANSFORMATIONS

LSC METHOD [RBF+ Const Poly]



```

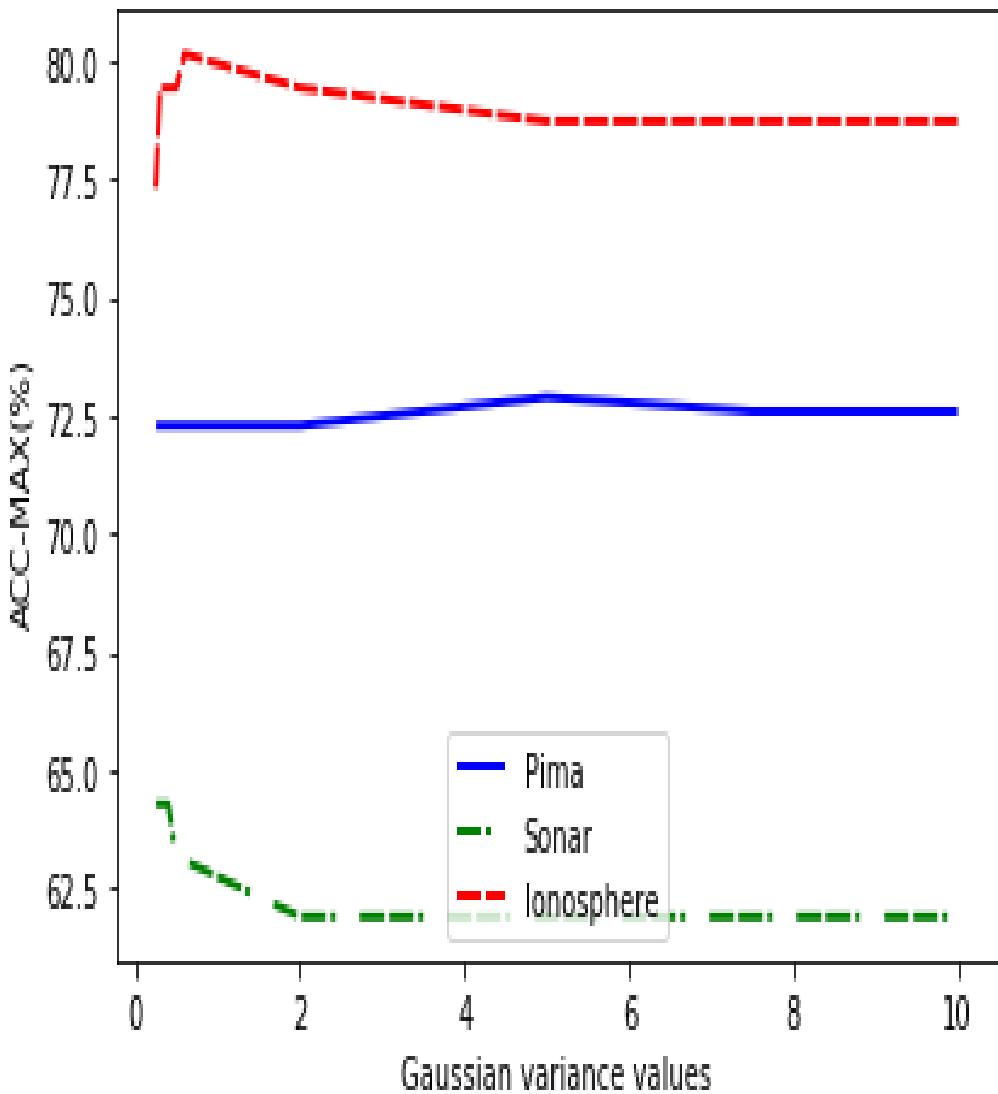
# /usr/bin/env python3
# -*- coding: utf-8 -*-
"""

Created on Sat Dec 1 19:18:27 2018
1
2 @author: sedna
3 ...
4 from IPython import get_ipython
5
6 get_ipython().magic('reset -sf')
7
8 import numpy as np
9 import pandas as pd
10 import statsmodels.api as sm
11 import statsmodels.discrete.discrete_model as smd
12 from statsmodels.stats.outliers_influence import variance_inflation_factor
13 from statsmodels.multivariate.pca import PCA as smPCA
14 from metrics_classifier import MetricClassifier
15 from sklearn import metrics
16
17 from sklearn.model_selection import train_test_split
18 from collections import OrderedDict, defaultdict
19 import matplotlib.pyplot as plt
20 from tools import Tools
21 from matplotlib import colors as mcolors
22 # colors for graphics with matplotlib and plotly
23 from smt.surrogate_models import rbf
24
25
26 testsize=0.4
27 #threshold=0.95
28 NDataSets=3
29
30 features_threshold=[95.0]
31 #features_threshold[90.0]
32 samples_transform=['sum','average','optim']
33 #samples_transform['sum']
34
35 Datasets=['Pima','Sonar','Ionosphere']
36 Datasets=['Ionosphere']
37
38 ac=MetricClassifier()
39
40 def show_best_pca_predictors(x,ncomp,threshold):
41     try:
42         eigenvalues=smPCA(x,ncomp).eigenvals
43     except:
44         eigenvalues=smPCA(x,ncomp,method='eig').eigenvals
45
46     order=[ii for ii,vals in sorted(enumerate(np.abs(eigenvalues)), \
47                                     key=lambda x:x[1].reverse=True)]
48     eigenvalues_PCA_sorted=[eigenvalues[order[ii]] for ii in range(ncomp)]
49     features_PCA_sorted=[x.columns[order[ii]] for ii in range(ncomp)]
50
51     key_colinearity= variance_influence_factors(x,features_PCA_sorted)
52
53     features_sorted=[]
54     eigenvalues_sorted=[]
55     for feat,eig in zip(features_PCA_sorted,eigenvalues_PCA_sorted):
56         if feat in key_colinearity():
57             if eig==key_colinearity()[0]:
58                 features_sorted.append(defeat)
59                 eigenvalues_sorted.append(eig)
60             else:
61                 features_sorted.append(feat)
62
63         eigenvalues_sorted.append(eig)
64
65
66     tot = sum(eigenvalues_sorted)
67     best_predictors=[]
68     cum_var_exp=0.0
69     for ii,lin in enumerate(eigenvalues_sorted[:]):
70         var_exp=(ii/tot)
71         cum_var_exp+=var_exp
72         if cum_var_exp < threshold:
73             best_predictors.append(features_sorted[ii])
74
75     return best_predictors
76
77 def variance_influence_factors(x,predictors,vif_threshold):
78
79

```

LSC METHOD [RBF + CONST POLY] COMPARISON BETWEEN DATASETS

ACCURACY BINARY CLASSIFICATION (TOP VALUES)



```
#!/usr/bin/env python3
# -*- coding: utf-8 -*-
"""
Created on Fri Nov  9 16:41:04 2018
@author: sedna
"""

import numpy as np
import pandas as pd
import statsmodels.api as sm
import statsmodels.discrete.discrete_model as smd
from statsmodels.multipivariate.pca import PCA as smPCA
from metrics_classifier import MetricClassifier
from sklearn import metrics
from sklearn.model_selection import train_test_split
from collections import OrderedDict
import matplotlib.pyplot as plt
from tools import Tools
from matplotlib import colors as mcolors
# colors for graphics with matplotlib and plotly
testsiz=0.4
treshold=0.90
NDataSets=3
Datasets=['Pima','Sonar','Ionosphere']
ac=MetricClassifier()

def show_best_pca_predictors(x,ncomp,treshold):
    try:
        eigenvalues=smPCA(x,ncomp).eigenvals
    except:
        eigenvalues=smPCA(x,ncomp,method='eig').eigenvals
    order=[ii for ii,vals in sorted(enumerate(np.abs(eigenvalues)),\n                                     key=lambda x:x[1],reverse=True)]
    eigenvalues_sorted=[eigenvalues[order[ii]] for ii in range(ncomp)]
    features_sorted=[x.columns[order[ii]] for ii in range(ncomp)]
    tot = sum(eigenvalues_sorted)
    best_predictors=[]
    cum_var_exp=0.0
    for ii in enumerate(eigenvalues_sorted):
        var_exp=(ii[0]/tot)
        cum_var_exp+=var_exp
        if cum_var_exp < treshold:
            best_predictors.append(features_sorted[ii])
    return best_predictors

def plot_matrix_matrix(X,Y,Title,xlabel,ylabel,Labels,
                      Linestyle,kind,scale,grid,text,boxstyle,mode):
    """
    Draw matrix vs matrix.
    kind of graphic
    0-plot
    1-scatter
    2-stem
    scale of plot
    'equal'
    'Log'
    'Semilogy'
    'Semilogx'
    'Loglog'
    mode of graphic
    0-full
    1-split
    """
    if (len(Y.columns)>7 and mode==0):
        colors = [lc for lc in mcolors.CSS4_COLORS.values() if lc !=(1,1,1)]
        keys = [kc for kc in mcolors.CSS4_COLORS.keys() if mcolors.CSS4_COLORS[kc] !=(1,1,1)]
```

PRELIMINARY CONCLUSIONS

The Lumped Sample Collocation Prediction Method is a global method which depends of Cumulative Distribution Function [cdf] used to compute the learning curve. Results seem to indicate that it has the same precision as linear methods. Best results were obtained adding features/sample columns and selecting the training points/labels as collocation points. Radial Basis Functions are employed to get the learning parameters.

<https://github.com/sednabcn/SiMLeng/tree/BinaryClassification-LSCollocation>

MACHINE LEARNING

SiMLeng-Scikit-learn package (in progress)

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 import pandas as pd
4 import numpy as np
5 import sklearn
6 import math
7 import nincreas
8 import seaborn as sns
9
10 from functools import wraps
11 from collections import OrderedDict
12 import matplotlib.pyplot as plt
13 from matplotlib import colors as mcolors
14 from bokeh.util import compat
15
16 from scipy import stats, linalg
17 from scipy.linalg import svd as skSVD
18 from scipy.stats import multivariate_normal
19 from mpl_toolkits.mplot3d import Axes3D
20
21 from tools import Tools
22
23 from sklearn import metrics
24 from sklearn import datasets
25 from sklearn import preprocessing
26 from sklearn import linear_model
27 from sklearn import semi_supervised
28
29 from sklearn.preprocessing import Binarizer,
30     FunctionTransformer, Imputer, KMeansCenterer,
31     LabelBinarizer, LabelEncoder, MultiLabelBinarizer,
32     MinMaxScaler, MaxAbsScaler, QuantileTransformer,
33     Normalizer, OneHotEncoder, PowerTransformer,
34     StandardScaler, add_dummy_feature,
35     PolynomialFeatures, binarize, normalize,
36     scale, robust_scale,
37     maxabs_scale, minmax_scale, label_binarize, quantile_transform
38
39 from nincreas import IR_on_increasing_size_pkl_2_svr
40
41 from sklearn.manifold import MDS as skMDS, smacof as skmacof, Isomap as sktsNE, LocallyLinearEmbedding as skLocalLEmbd, spectral_embedding as skSpectEmbed
42
43 from sklearn.model_selection import cross_val_score
44 from sklearn.model_selection import train_test_split
45 from sklearn.decomposition import PCA as skPCA
46
47 from sklearn.grid_search import GridSearchCV
48 from sklearn.model_selection import KFold
49
50
51 from sklearn.feature_selection import SelectKBest as skSB
52 from sklearn.feature_selection import SelectFValueFilter as skFVF
53 from sklearn.feature_selection import LassoLarsCV as skLassoLarsCV
54 from sklearn.feature_selection import VarianceThreshold as skVT
55 from sklearn.feature_selection import RFE as skRFE
56 from sklearn.feature_selection import RFECV as skRFECV
57 from sklearn.feature_selection import RFECV as skRFECV
58 from sklearn.feature_selection import SelectFromModel as skSFM
59 from sklearn.feature_selection import SelectFromModel as skFMD
60 from sklearn.grid_search import GridSearchCV as skGridSCV
61
62 from sklearn.learning_curve import learning_curve as skLearnCurve, validate_curve as skValCurve
63
64 from sklearn.pipeline import make_pipeline
65
66 from sklearn.pipeline import Pipeline
67
68
69 from sklearn import linear_model as sklm
70 from sklearn.linear_model import ARDRegression, BayesianRidge, ElasticNet, ElasticNetCV, Hinge, Huber, HuberRegressor, Lars, LarsCV, Lasso, LassoCV, LassoLars, LassoLarsCV, LassoLarsIC, Line
71 Perception, RANSACRegressor, RandomizedLasso, RandomizedLogisticRegression, Ridge, RidgeCV, \
72 RidgeClassifier, RidgeClassifierCV, SGDClassifier, SquaredLoss, TheilSenRegressor
73
74 from sklearn import linear_model as sklm
75 from sklearn import base, bayes, cd_fast, coordinate_descent, \
76 enet_path, huber, lars_path, lasso_stability, path_coordinate_descent, \
77 logistic, logreg, orthogonal_mp, orthogonal_mp_gram, passive_aggressive, perceptron, \
78 randomized_l1, renyi, ridge, ridge_regression, sag, sagd, fast, sgd, fast_mcmc, gradient, \
79 theil_sen
80
81 from sklearn import neighbors as skn
82
83 from sklearn import random_projection as skrp
84 from sklearn import manifold as skm
85
86 from sklearn import svm
87 from sklearn.svm import SVC, NuSVC, SVR, NuSVR, OneClassSVM, LinearSVC, \
88 LinearSVR
89
90 from sklearn import svm as skSVC
91 from sklearn import linear_model as sklsvr
92 from sklearn import isotonic as skir
93 from sklearn import tree as sktr
94 from sklearn import ensemble as sken
95
96 from sklearn import neural_network as sknn
97
98 from sklearn import gaussian_process as skgpc
99 from sklearn import gaussian_process as skgpr
100 from sklearn import gaussian_process as skgpp
101 from sklearn import gaussian_process as skgpc
102 from sklearn import gaussian_process.correlation_models as skgc
103
104 from sklearn import gaussian_process.kernels as skkern
105 from sklearn import gaussian_process.kernels as skkern
106 from sklearn import gaussian_process.kernels as skkern
107
108 from sklearn import kernel_approximation as skkappa
109 from sklearn import kernel_approximation as skkappa
110 from sklearn import kernel_approximation as skkappa
111
112 from sklearn import kernel as skker
113 from sklearn import kernel as skker
114
115 from sklearn import kernel_ridge as skkr
116 from sklearn import kernel_ridge as skkr
117
118 from sklearn import discriminant_analysis as skda
```

MACHINE LEARNING

Scikit-learn package

DECISION TREE

Banknote Authentication Dataset / Iris Dataset

```
Python 3.7.3 (default, Apr 3 2019, 05:39:12)
Type "copyright", "credits" or "license" for more information.
IPython 7.5.0 -- An enhanced Interactive Python.
Restaring kernel...
```

```
runfile('/home/sedna/PYTHON/PYML/DECISION-TREE/Decision-Tree-01.py', wdir='/home/sedna/PYTHON')
Dataset Length: 400
Dataset: 0 1 2 3 4
0 R 1 1 1 2
1 R 1 1 1 3
2 R 1 1 1 4
3 R 1 1 1 5
4 R 1 1 1 6
Results Using Gini Index:
```

```
Predicted values:
```

```
[R: 188 L: 212]
[True: 188 False: 212]
[0: 150 1: 50 2: 50]
Confusion Matrix: [[0 6 7]
[0 19 71]
[0 19 71]]
Accuracy : 73.40425531914693
```

```
Report:
```

```
          precision    recall   f1-score   support
R       0.00      0.00      0.00     12
L       0.73      0.79      0.76     85
R       0.74      0.79      0.76     90
In [1]: rec = 0.74
        precision = 0.74
        recall = 0.74
        f1-score = 0.74
        support = 188
        weighted avg = 0.69 0.51 0.51 188
Results Using Entropy:
```

```
Predicted values:
[True: 188 False: 212]
[0: 100 1: 49 2: 51]
Confusion Matrix: [[0 6 7]
[0 20 79]
[0 20 79]]
Accuracy : 70.74466085106383
```

```
Report:
```

```
          precision    recall   f1-score   support
R       0.00      0.00      0.00     12
L       0.71      0.74      0.72     85
R       0.71      0.78      0.74     90
accuracy = 0.74
precision = 0.74
recall = 0.74
f1-score = 0.74
weighted avg = 0.69 0.51 0.49 188

```

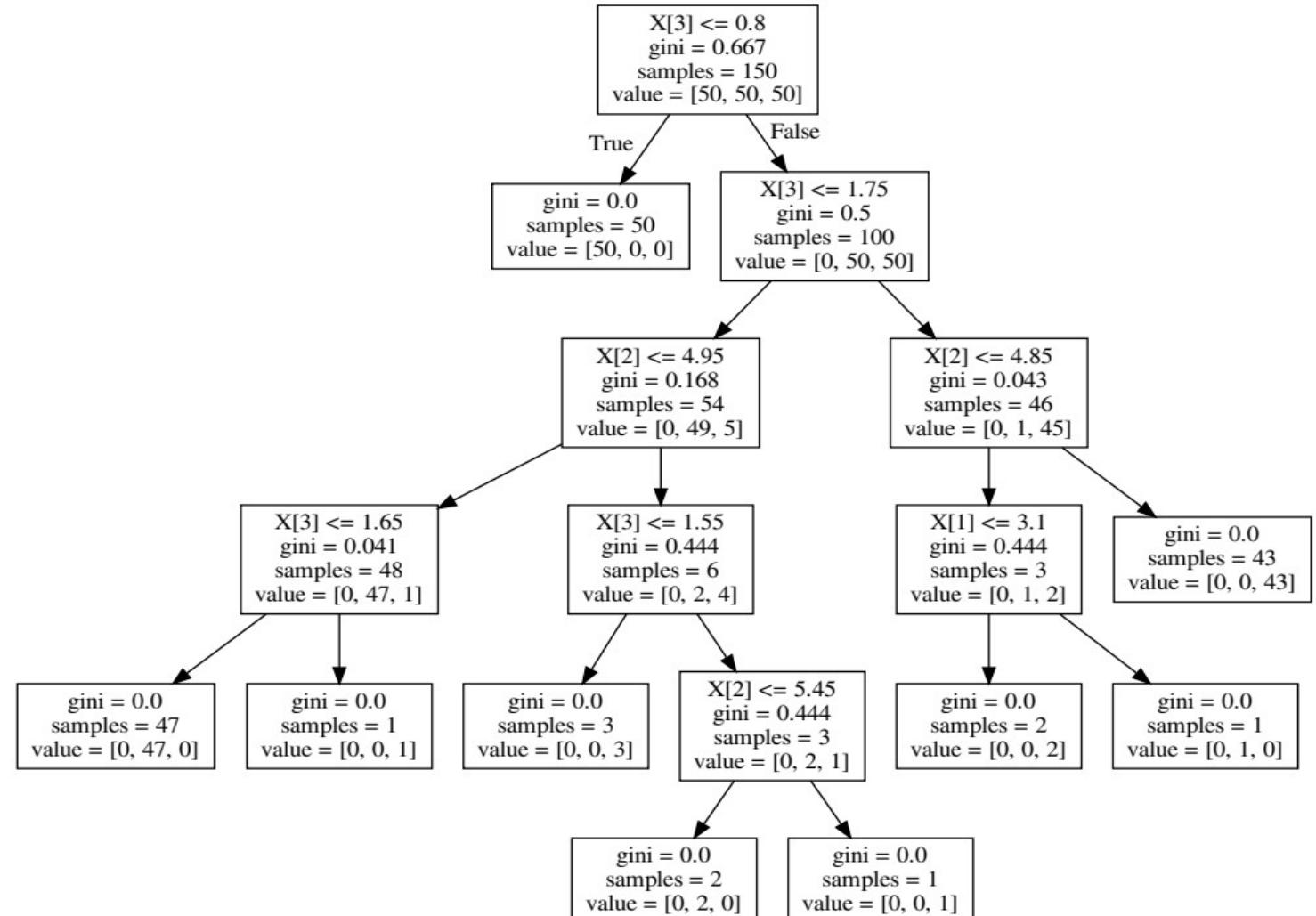
```
runfile('/home/sedna/PYTHON/PYML/DECISION-TREE/Decision-Tree-02.py', wdir='/home/sedna/PYTHON')
no predicted samples
[precision, recall, f1-score, support]
no predicted samples
[precision, recall, f1-score, support]
```

```
In [2]: runfile('/home/sedna/PYTHON/PYML/DECISION-TREE/Decision-Tree-02.py', wdir='/home/sedna/PYTHON')
          2.77154716 1.7847839599 0
0 1.78571 1.169761 0
1 3.678320 2.812814 0
2 3.989209 2.209014 0
3 2.999209 2.209014 0
4 2.999209 2.209014 0
[1 1.2857131 1.16976141 0, 1
[3.67831985 2.81281357 0, 1
[3.9892092 2.20901421 0, 1
[2.9992092 2.20901421 0, 1
[9.00220326 3.33964719 1, 1
[7.44454233 0.47668338 1, 1
[10.00000000 0.00000000 1, 1
[6.64228735 3.33998376 1, 1]
Sprint = 6.642
[X1 < 1.729]
[X1 < 1.729]
[0,0]
[X1 < 7.498]
[X1 < 7.495]
[1,0]
[1,0]
[X1 < 7.498]
[1,0]
Expected=0, Got=0
Expected=0, Got=0
Expected=0, Got=0
Expected=1, Got=1
Expected=1, Got=1
Expected=1, Got=1
Expected=1, Got=1
Expected=1, Got=1

```

```
In [3]: runfile('/home/sedna/PYTHON/PYML/DECISION-TREE/Decision-Tree-03.py', wdir='/home/sedna/PYTHON')
          0 1 2 3 4
0 3.62160 -8.6661 -2.8073 -0.4460 0
1 4.32060 -2.6380 1.9242 0.46210 0
2 3.86600 -2.6380 1.9242 0.10645 0
3 3.86600 9.3535 0.0116 -0.58840 0
4 0.329560 4.8552 4.0116 -0.58840 0

```



MACHINE LEARNING

Scikit-learn package

kMeans

An Artificial Dataset

```
Python 3.7.3 (default, Apr 3 2019, 05:39:12)
Type "copyright", "credits" or "license" for more information.
```

```
[Python 3.7.3] In [1]: As an interactive Python
notebook[home/edua/PYTHON/PYML/kMeans/kMeans.ipynb, vdr:/home/edua/PYTHON/PYML/kMeans]
```

```
[380, 3]
[[3, 40],
 [19, 40],
 [18, 53],]]
```

```
[1]
```

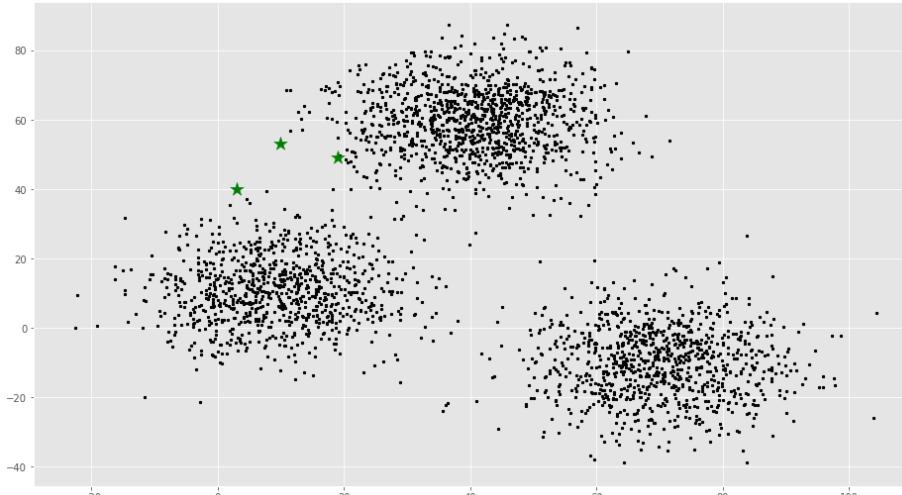
```
[2]
```

```
In [2]: runfile('/home/edua/PYTHON/PYML/kMeans/kMeans.ipynb', vdr='/home/edua/PYTHON/PYML/kMeans')
[380, 3]
[[380.0098 0.7650000,
 [3.809207 0.7699882,
 [78.533674 0.1036521],
 [61.679025 0.1070829],
 [40.0540139 0.8484874],]]
```

```
[3]
```

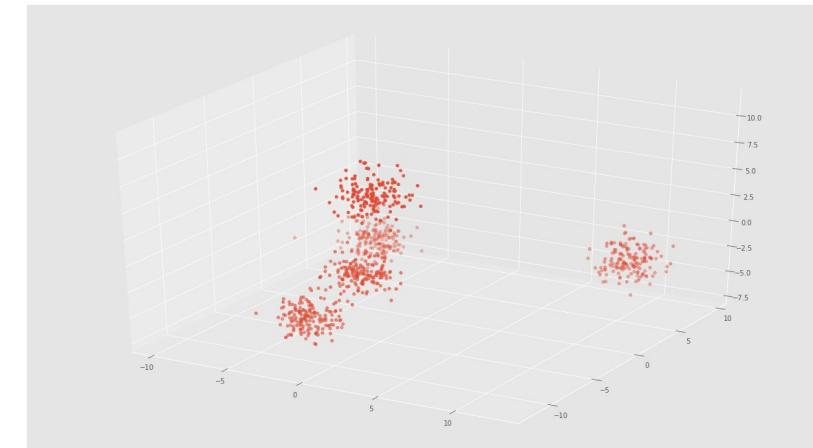
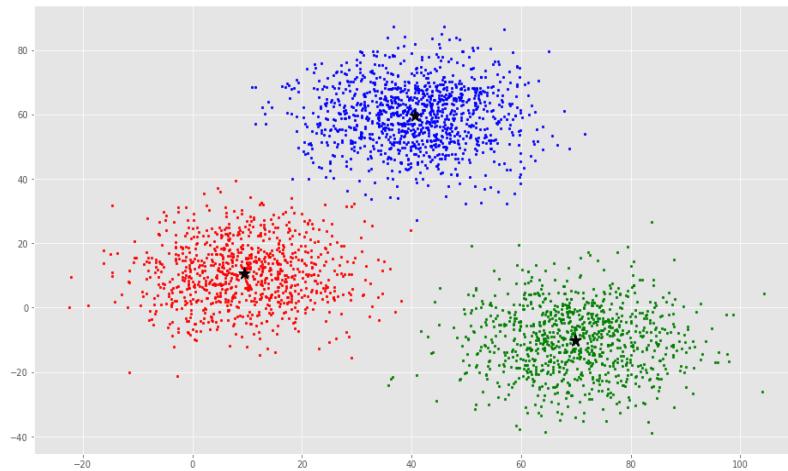
```
[4]
```

```
In [3]:
```



k=3

k=4



MACHINE LEARNING

Scikit-learn package

k-Nearest-Neighbor

Iris Dataset

ACCURACY

Train Set	Test Set	k=3	k=5	k=10
66%	34%	60.4878048780487%	62.926829268292686%	65.3658536585366%
80%	20%	62.3255813953488%	62.32558139534884%	68.83720930232559%

MACHINE LEARNING

Scikit-learn package

Naive-Bayes

ACCURACY

Train Set	Test Set	Pima-Indians Dataset	Iris Dataset	Ionosphere Dataset
67%	33%	74.0157480314960%	98.0%	80.17241379310344%
80%	20%	74.6753246753246%	90.0%	88.73239436619718%

MACHINE LEARNING

Scikit-learn package

Random-Forest

- Petrol-consumption / Breast-cancer-wisconsin Dataset

```
Python 3.7.3 (default, Apr 3 2019, 05:39:12)
Type "copyright", "credits" or "license" for more information.
```

```
IPython 7.5.0 - An enhanced Interactive Python
runfile('/home/sedna/PYTHON/PYML/RANDOM-FOREST/forest-40.py', wdir='/home/sedna/PYTHON/PYML/RANDOM-FOREST')
Petrol_tax Average_income ... Population_Driver_licence(%) Petrol_Consumption
```

0	9.0	3571	0.525	541
1	9.0	4092	0.572	524
2	9.0	3865	0.580	561
3	7.5	4870	0.529	414
4	8.0	4399	0.544	410

```
[5 rows x 5 columns]
Mean Absolute Error: 51.76500000000001
Mean Squared Error: 4216.166749999999
Root Mean Squared Error: 64.93201637097064
```

```
In [2]: runfile('/home/sedna/PYTHON/PYML/RANDOM-FOREST/forest-0clasiif.py', wdir='/home/sedna/PYTHON/PYML/RANDOM-FOREST')
Trees: 1
Scores: [56.09756097560976, 63.414631463146, 60.97560975609756, 58.536585365853654, 73.17073170731707]
Mean Accuracy: 62.43%
Trees: 5
Scores: [70.73170731707317, 58.536585365853654, 85.36585365853658, 75.60975609756098, 63.414631463146]
Mean Accuracy: 70.732%
Trees: 10
Scores: [75.60975609756098, 80.48780487804879, 92.6829268292683, 73.17073170731707, 70.73170731707317]
Mean Accuracy: 78.53%
```

```
In [3]: runfile('/home/sedna/PYTHON/PYML/RANDOM-FOREST/forest-001.py', wdir='/home/sedna/PYTHON/PYML/RANDOM-FOREST')
CodeNumber ClumpThickness ... Minoses CancerType
count 6.800000e+02 698.000000 ... 698.000000 698.000000
mean 1.071807e-06 4.41695 ... 1.59238 2.69054
std 6.175324e-05 2.817673 ... 1.716162 0.951596
min 6.163400e-04 1.000000 ... 1.000000 2.000000
25% 8.702582e-05 2.000000 ... 1.000000 2.000000
50% 1.171710e-06 4.000000 ... 1.000000 2.000000
75% 1.238354e-06 6.000000 ... 1.000000 4.000000
max 1.345454e-07 10.000000 ... 10.000000 4.000000
```

```
[8 rows x 10 columns]
Train_x Shape :: (477, 9)
Train_y Shape :: (477,)
Test_x Shape :: (205, 9)
Test_y Shape :: (205,)
Trained model :: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
max_depth=None, max_features='auto', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=10,
n_jobs=None, oob_score=False, random_state=None,
verbose=0, warm_start=False)
Actual outcome :: 2 and Predicted outcome :: 2
Actual outcome :: 2 and Predicted outcome :: 2
Actual outcome :: 2 and Predicted outcome :: 2
Actual outcome :: 4 and Predicted outcome :: 4
Actual outcome :: 4 and Predicted outcome :: 4
Train Accuracy :: 1.0
Test Accuracy :: 0.9609756097560975
Confusion matrix [[126 5]
 [ 3 71]]
/usr/local/lib/python3.7/dist-packages/scikit_learn-0.21.2-py3.7-linux-x86_64.egg/sklearn/ensemble/forest.py:245: FutureWarning: The default value of n_estimators will change from 10 in version 0.20 to 100 in 0.22.
"10 in version 0.20 to 100 in 0.22.", FutureWarning)
```

Petrol-consumption Dataset
TREES: 10
MEAN ACCURACY: 78.537%

Breast-cancer-wisconsin Dataset
TREES: 10
MEAN ACCURACY: 96.09756%

MACHINE LEARNING

Scikit-learn package

SVM LINEAR CLASSIFICATION (TEST SIZE=30%)

METRICS	Breath-cancer dataset	Iris dataset	Wine dataset	Digits dataset (class =10)
ACCURACY	98.2456140%	93.3333333%	98.1481481%	96.6666666%
PRECISION	98.2469120%	93.4320987%	98.2322323%	96.6742587%
RECALL	98.2456140%	93.3333333%	98.1481481%	96.6666666%

SVM LINEAR PENALIZED 'l1' CLASSIFICATION (TEST SIZE=30%)

METRICS	Breath-cancer dataset	Iris dataset	Wine dataset	Digits dataset (class =10)
ACCURACY	98.8304093%	93.3333333%	100.0%	96.1111111%
PRECISION	98.8516746%	93.4558823%	100.0%	96.1411426%
RECALL	98.8304093%	93.3333333%	100.0%	96.1111111%

Applied AI and Deep Learning

Build and Serialisation Keras Model

Classification using Reuters Dataset

PARAMETERS

X_train_rows	8982
X_test_rows	2246
max_words	1000
num_classes	46
test_split	0.2
validation_split	0.1
Epochs	3
batch_size	32
Optimizer	Adam

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 512)	512512
dropout (Dropout)	(None, 512)	0
dense_1 (Dense)	(None, 46)	23598
<hr/>		

Total params: 536,110

Trainable params: 536,110

Non-trainable params: 0

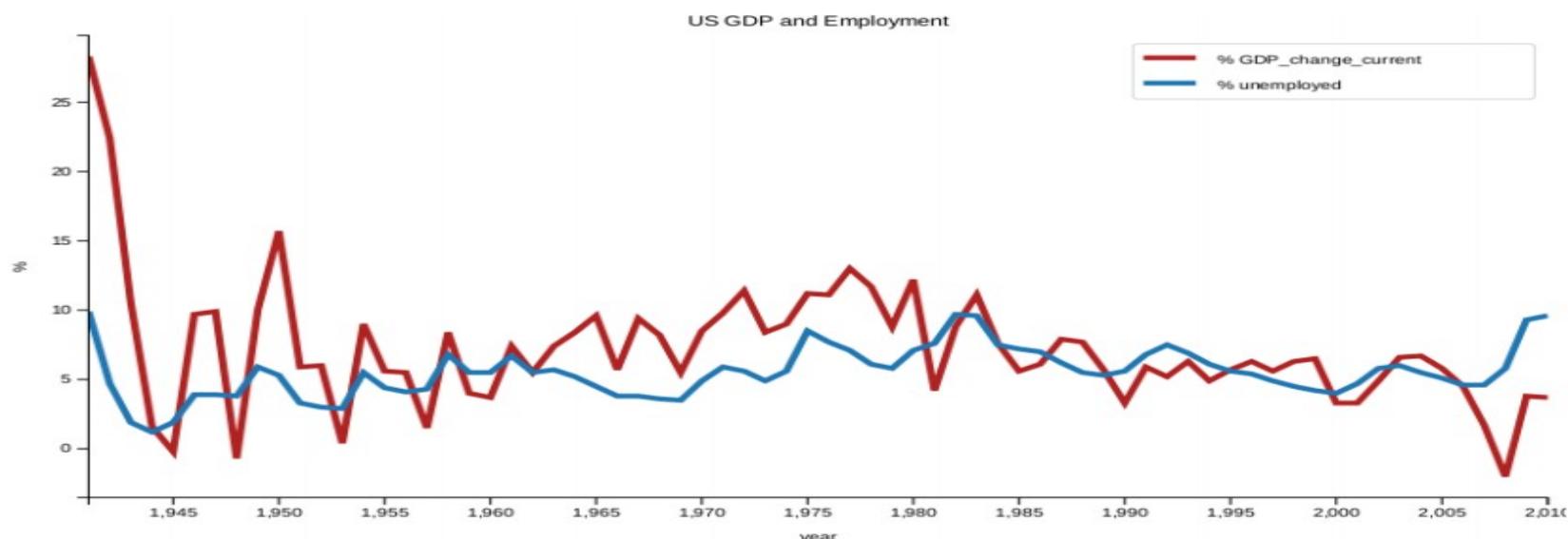
HISTORY

Epoch	val_loss	val_acc	loss	acc
0	0.969766	0.787542	1.440816	0.680193
1	0.851318	0.810901	0.783447	0.820859

Mpld3 Project

Analyzing US Economic Data and Building a Dashboard

Gross domestic product (GDP) is a measure of the market value of all the final goods and services produced in a period, often annually. GDP is an indicator of how well the economy is doing. A drop in GDP indicates the economy is producing less; similarly an increase in GDP suggests the economy is performing better. GDP is considered the "world's most powerful statistical indicator of national development and progress".[1] GDP can be determined in three ways, all of which should, in principle, give the same result. They are the production (or output or value added) approach, the income approach, or the speculated expenditure approach. Literature 1.-Lepenies, Philipp (2016). The Power of a Single Number: A Political History of GDP. New York: Columbia University Press. ISBN 978-0-231-17510-4.



Applied AI and Deep Learning

Build and Serialisation Keras Model

Classification using Pimas-Indians-Diabetes Dataset

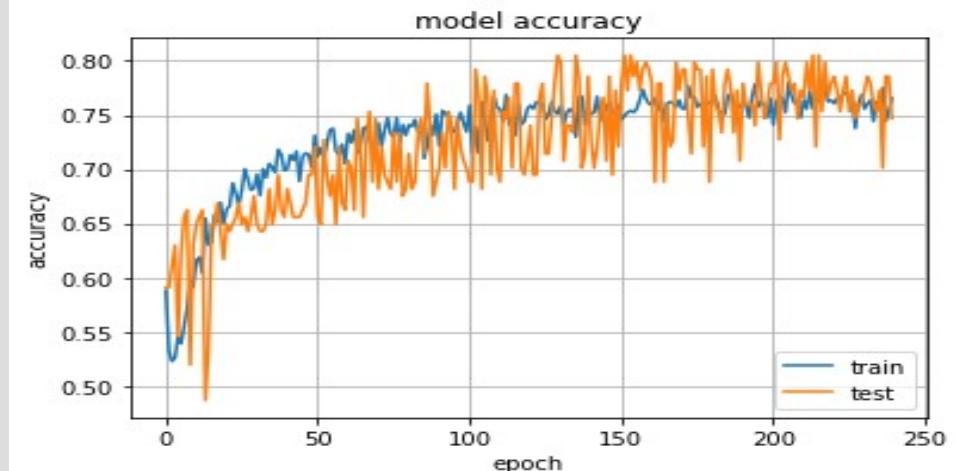
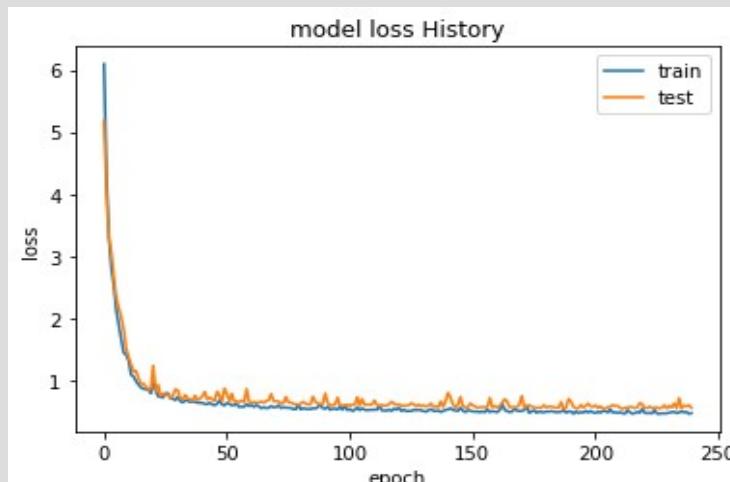
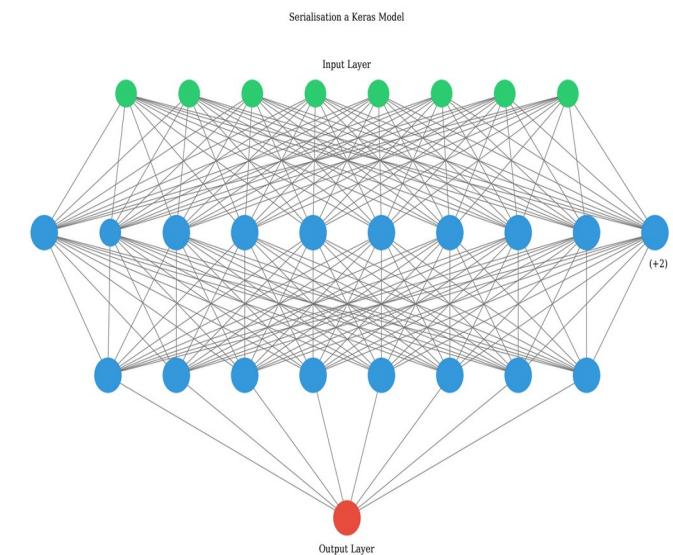
PARAMETERS

X_rows	768
features_size	8
num_classes	1
validation_split	0.2
Epochs	250
batch_size	32
Optimizer	Adamax

Model: "sequential"

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 12)	108
dense_2 (Dense)	(None, 8)	104
dense_3 (Dense)	(None, 1)	9

Total params: 221
Trainable params: 221
Non-trainable params: 0



Anomaly Detection case using IBM Watson Study platform

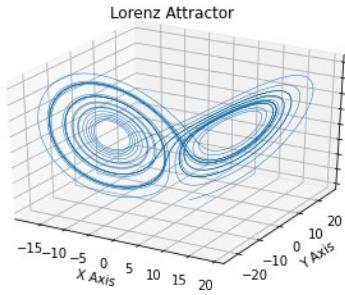
Taking data from <https://raw.githubusercontent.com/romeokienzler/developerWorks/master/lorenzattractor/>

We have implemented re-use the unsupervised anomaly detection algorithm but turn it into a simpler feed forward neural network for supervised classification. For this purpose, we have trained the neural network from healthy and broken samples and at later stage hook it up to a message queue for real-time anomaly detection.

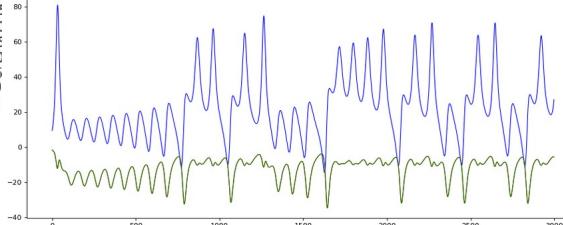
Model: "sequential"

Layer (type)	Output Shape	Param
--------------	--------------	-------

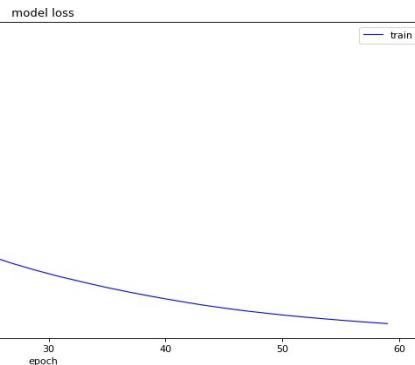
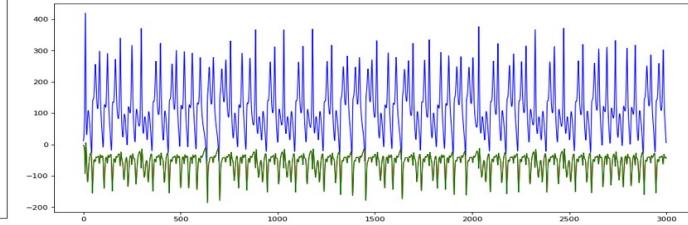
dense_1 (Dense)
(None, 3000)
9003000



Healthy frequencies



Broken frequencies



Healthy frequencies
Score

0.92995715

0.9455121

Broken frequencies
Score

0.07092416

0.05365459

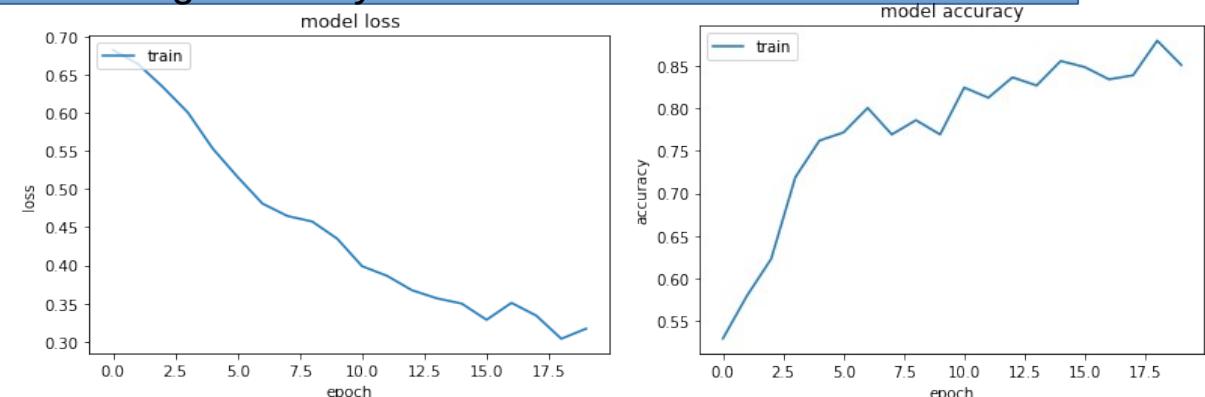
A KERAS MODEL RUNNING IN DL4J ON SPARK

We have implemented a NN for Sonar Dataset. The file "sonar.mines" contains 111 patterns obtained by bouncing sonar signals off a metal cylinder at various angles and under various conditions. The file "sonar.rocks" contains 97 patterns obtained from rocks under similar conditions. The transmitted sonar signal is a frequency-modulated chirp, rising in frequency. The data set contains signals obtained from a variety of different aspect angles, spanning 90 degrees for the cylinder and 180 degrees for the rock. Each pattern is a set of 60 numbers in the range 0.0 to 1.0. Each number represents the energy within a particular frequency band, integrated over a certain period of time. The integration aperture for higher frequencies occur later in time, since these frequencies are transmitted later during the chirp. The label associated with each record contains the letter "R" if the object is a rock and "M" if it is a mine (metal cylinder). The numbers in the labels are in increasing order of aspect angle, but they do not encode the angle directly.

Model: "sequential" Layer (type) Output
Shape Param

dense (Dense) (None, 60) 3660
dense_1 (Dense) (None, 30) 1830
dense_2 (Dense) (None, 2) 62

Total params: 5,552 Trainable
params: 5,552



```
!$SPARK_HOME/bin/spark-submit \
--class skymind.dsx.KerasImportCSVSparkRunner \
--files sonar.csv,my_modelx.h5 \
--master $MASTER \
dl4j-snapshot.jar \
-batchSizePerWorker 15 \
-indexLabel 60 \
-train false \
-numClasses 2 \
-modelFileName my_modelx.h5 \
-dataFileName sonar.csv
```

SCORES

# of classes:	2
Accuracy:	0.8894
Precision:	0.8894
Recall:	0.8912
F1 Score:	0.8856

SiMLeng PROJECT

**Tourism enhancers index
by neighborhoods in
London**

RUPERTO P. BONET

rupertobonet@gmail.com

DESCRIPTION OF THE PROJECT

London is one of the most popular tourist destination in the whole world, join to Paris and Bangkok. Each and every year London city attracts millions of visitors from other countries. During 2019 London received 19.09 millions of tourists according to the Mastercard's annual Global Cities Index publishing.

The tourism industry is one of the engines of the economy growth in UK. London first works int the Tourism Sector Deal proposal. This proposal defends and enhances the Lodon's role as a gateway , open and welcoming destination for both business and leisure visitors

[<https://www.londonfirst.co.uk/what-we-do/economy-and-tax/tourism>]

In this project, motivated by the challenge of understanding London urban environments, and based on the opportunities created by geoenabled social data, we address the problem of comparing boroughs and neighborhoods in London city .The problem we study has applications to recommending locations in London city as tourist enhancers. Imagine an entrepreneur decides by looking for areas with tourist profile to open a new establishment related to the tourist target like as in the neighbourhoods tradillionaly most visited, for example, where should I open a basement to rent bycicles?: the methods developed in this paper allow to match each neighborhood in a borough with the most similar neighborhood in an other borough, or any other borough that someone wishes to compare.

Tourist-enhancers-index is a mean value index of turist activities by neighbourhood. This activities include the following categories: arts_entertainment, building, education, event, food, nightlife, parks_outdoors, shops and travel. Data are collected of several sources including Foursquare api.

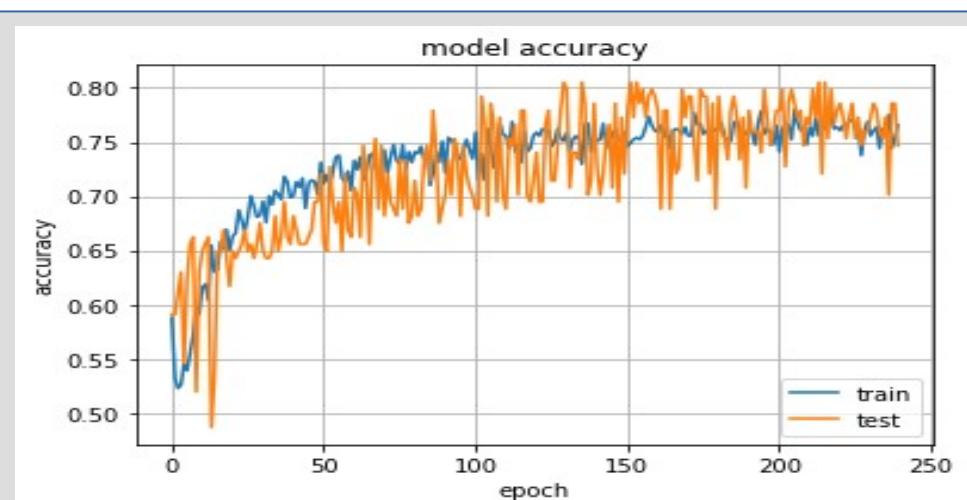
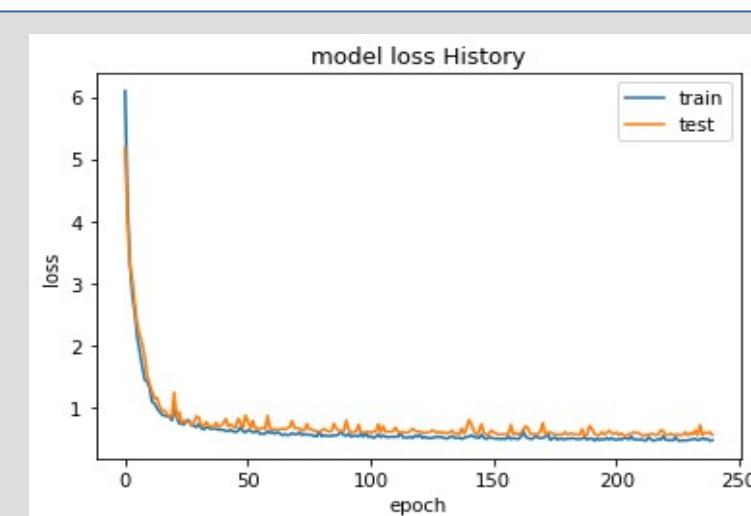
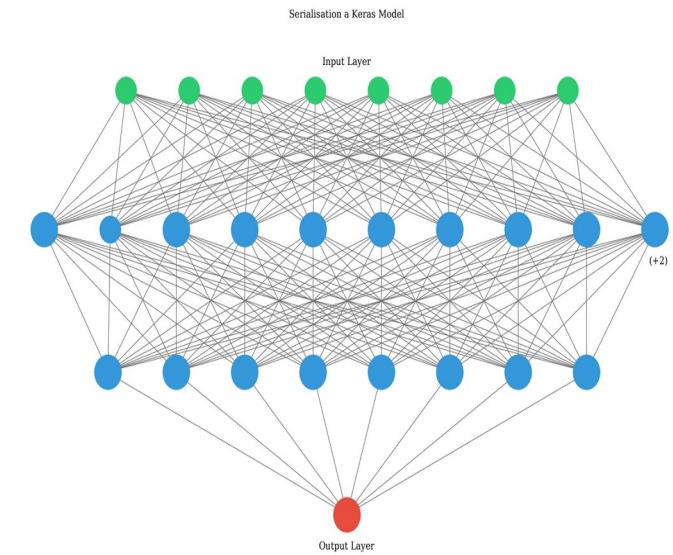
Applied AI and Deep Learning

Build and Serialisation Keras Model

Classification using Pimas-Indians-Diabetes Dataset

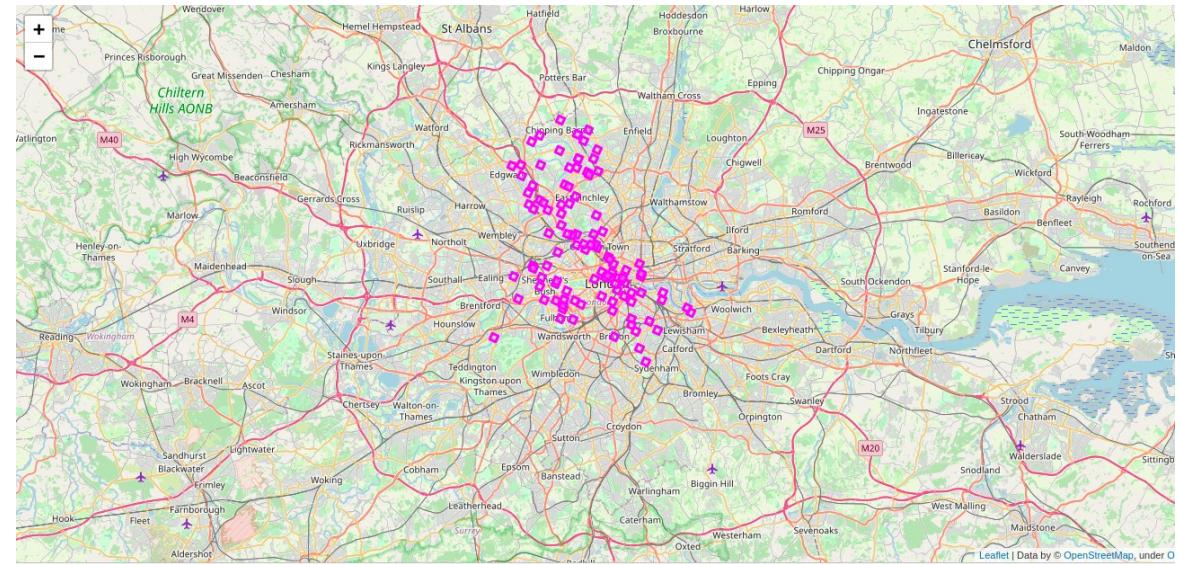
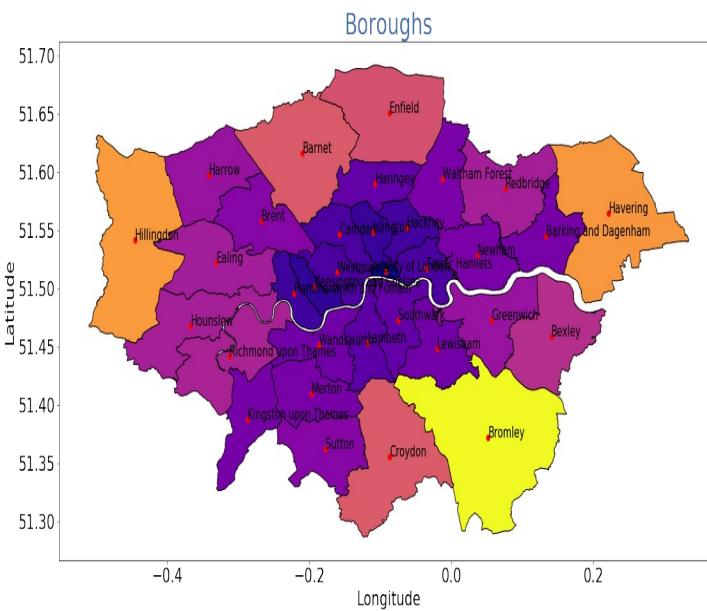
PARAMETERS	
X_rows	768
features_size	8
num_classes	1
validation_split	0.2
Epochs	250
batch_size	32
Optimizer	Adamax

```
Model: "sequential"
-----  
Layer (type)      Output Shape       Param #  
-----  
dense_1 (Dense)   (None, 12)        108  
-----  
dense_2 (Dense)   (None, 8)         104  
-----  
dense_3 (Dense)   (None, 1)         9  
-----  
Total params: 221  
Trainable params: 221  
Non-trainable params: 0
```

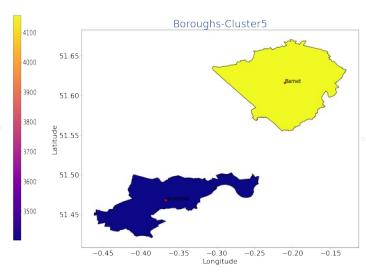
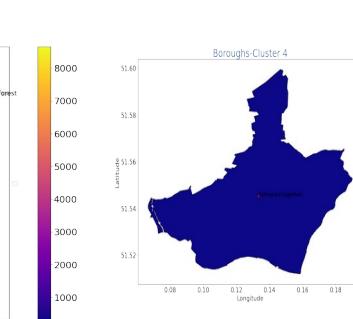
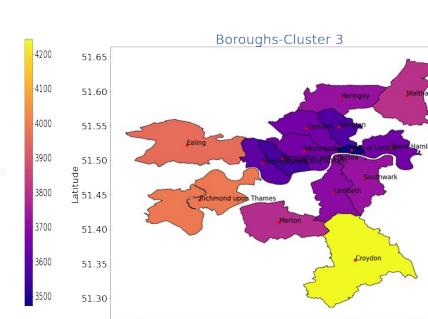
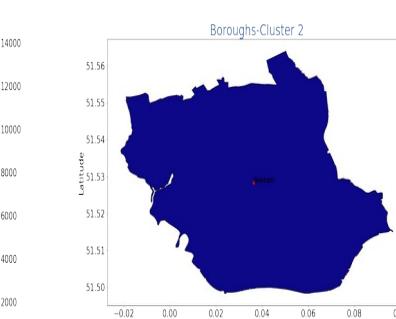
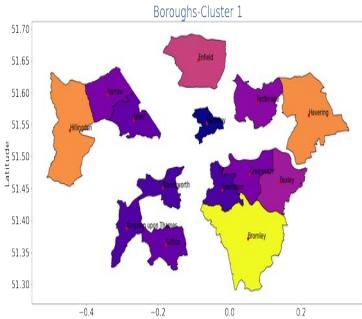


TOURIST-ENHANCERS INDEX BY NEIGHBOURHOODS IN LONDON

MOST VISITED PLACES IN LONDON



KMEANS CLUSTERING (K=5)



TOURIST-ENHANCERS INDEX BY NEIGHBOURHOODS IN LONDON

RECOMENDATION ENGINE GOOD CLASSIFIER (FAVOURITE-NON FAVOURITE)

PARAMETERS

X_rows	515
features_size	9
num_classes	1
validation_split	0.25

SUPERVISED CLASSIFIERS ACCURACY

svm: 0.7984496124031008

KNN: 0.7906976744186046

Logistic

Regresion: 0.7906976744186046

Decision

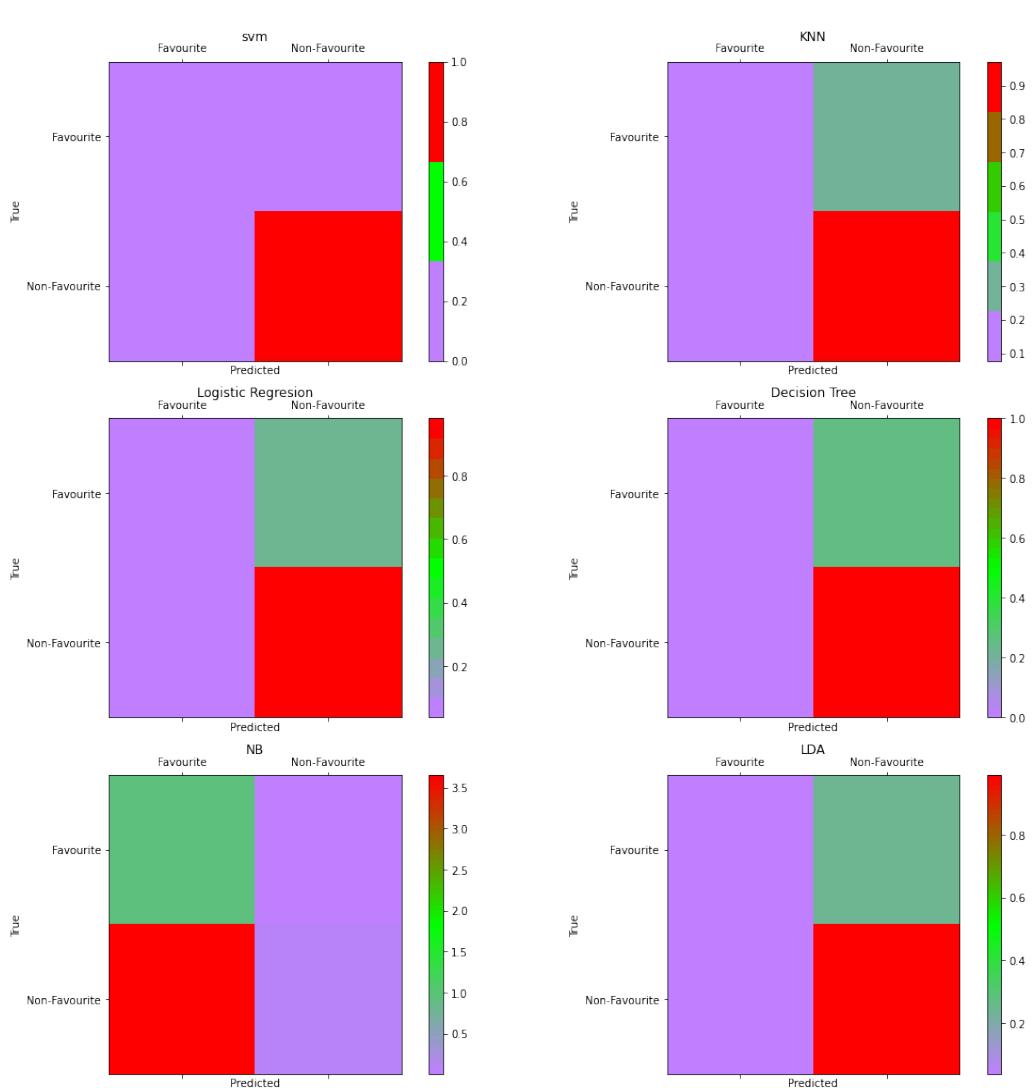
Tree: 0.7984496124031008

NB: 0.2558139534883721

LDA: 0.7984496124031008

CONFUSION MATRIX

Confusion matrix of Classifier Models



TOURIST-ENHANCERS INDEX BY NEIGHBOURHOODS IN LONDON

RECOMENDATION ENGINE

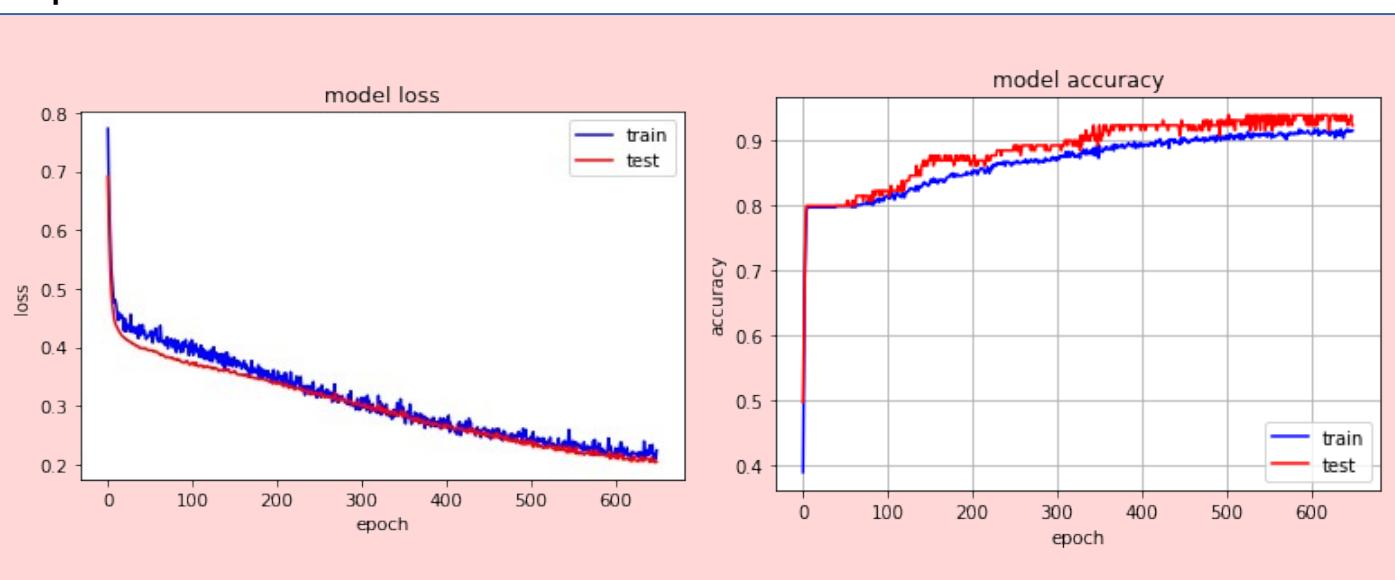
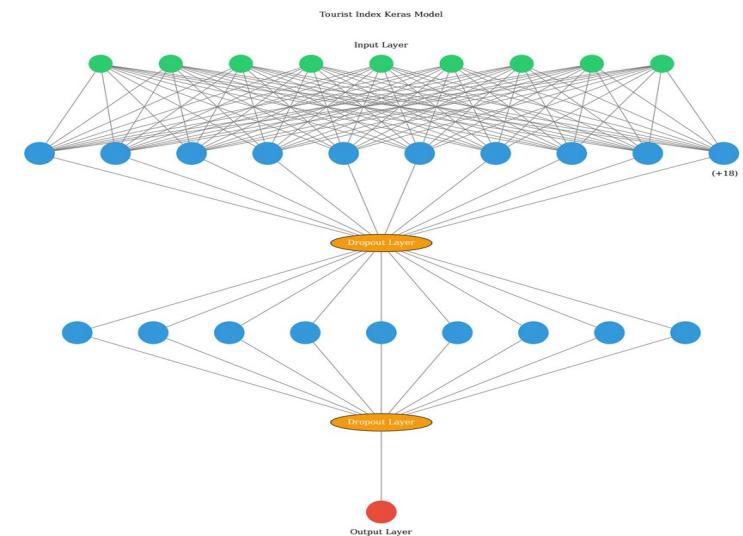
GOOD CLASSIFIER (FAVOURITE-NON FAVOURITE)

PARAMETERS

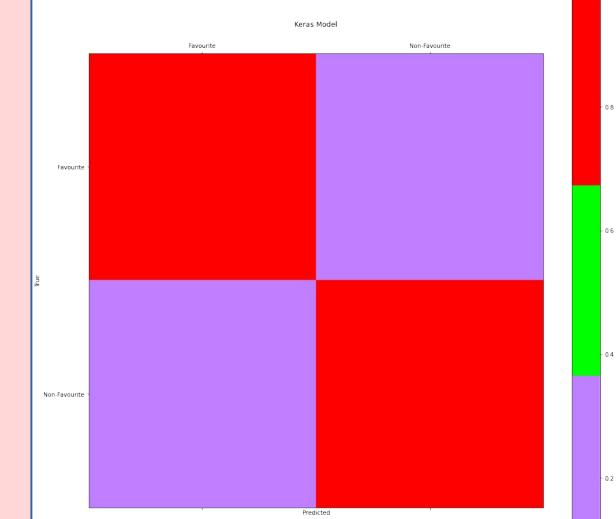
X_rows	515
features_size	9
num_classes	1
validation_split	0.25
Epochs	250
batch_size	50
Optimizer	Adam

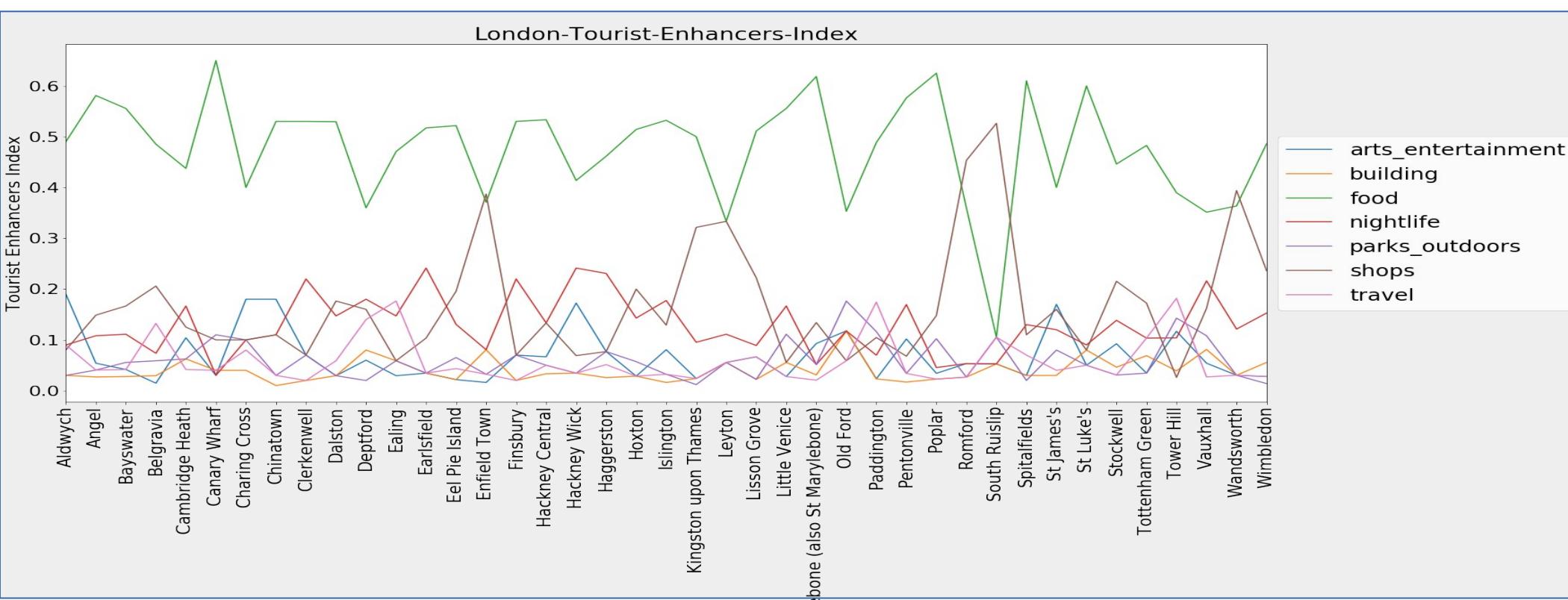
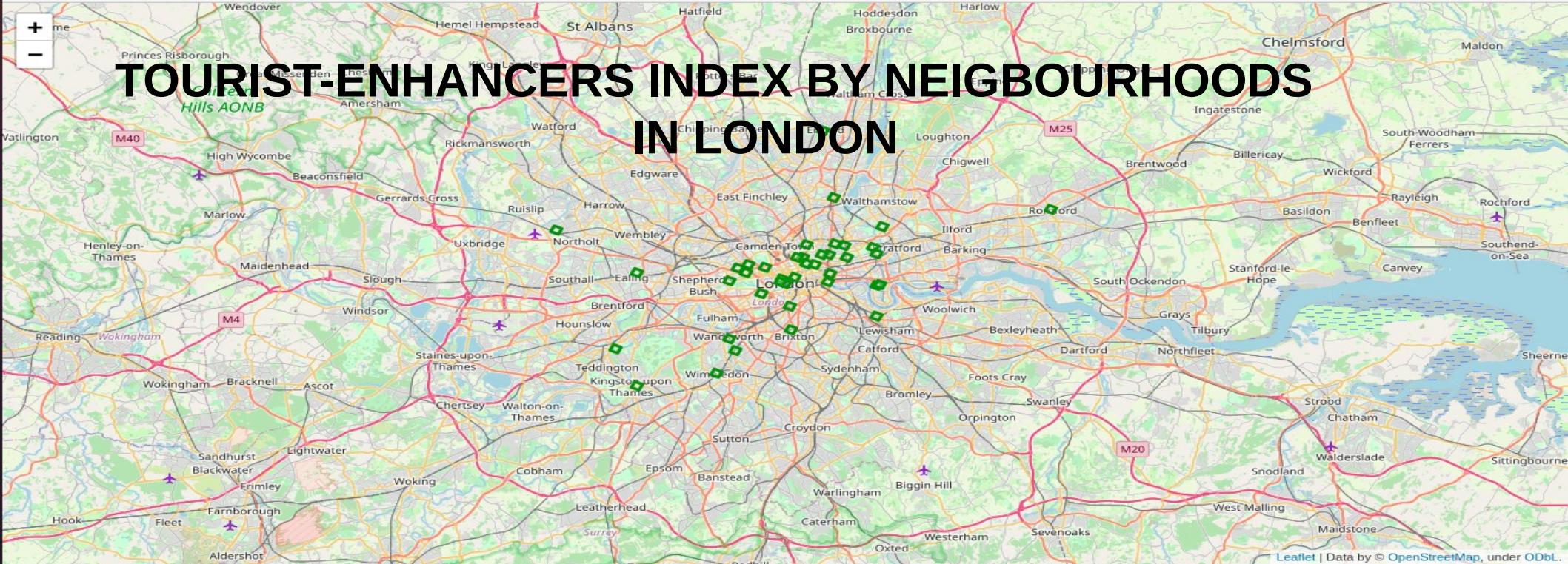
Model: "sequential"

Layer (type)	Output Shape	Param #
dense_6 (Dense)	(None, 28)	280
dropout_2 (Dropout)	(None, 28)	0
dense_7 (Dense)	(None, 9)	261
dropout_3 (Dropout)	(None, 9)	0
dense_8 (Dense)	(None, 1)	10
<hr/>		
Total params: 551		
Trainable params: 551		
Non-trainable params: 0		



CONFUSION MATRIX





Applied AI and IMAGE PROCESSING

KERAS : MNIST HANDWRITTEN DIGIT RECOGNITION

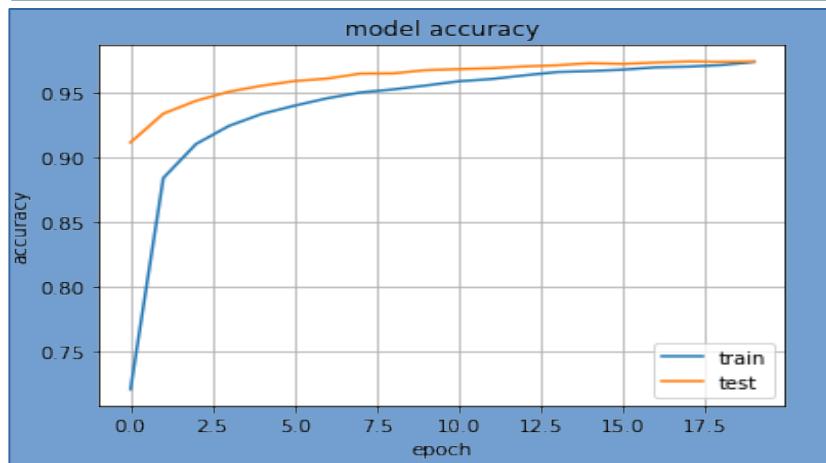
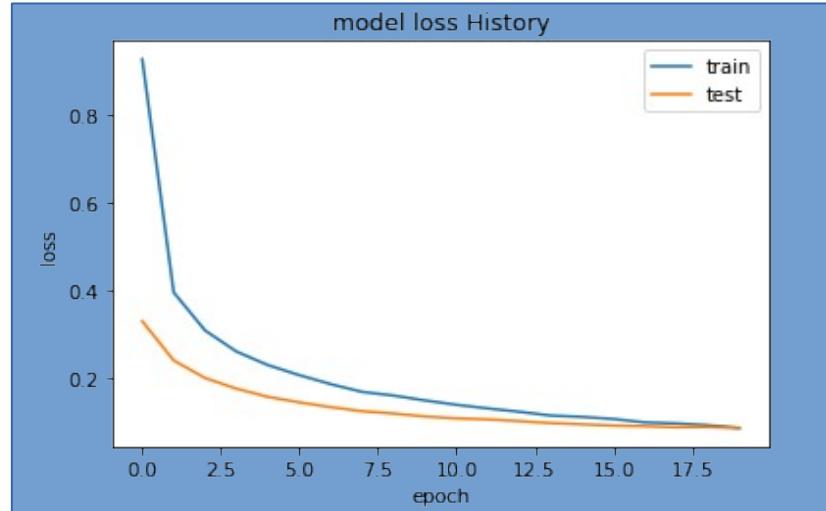
PARAMETERS

train samples 60000
test samples 10000
classes 10
hidden layers 128
dropout 0.25
learning rate 0.0002
batch_size 128
Validation split 0.2

Model: "sequential_1"

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 128)	100480
dropout_1 (Dropout)	(None, 128)	0
dense_2 (Dense)	(None, 128)	16512
dropout_2 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 10)	1290

Total params: 118,282
Trainable params: 118,282
Non-trainable params: 0



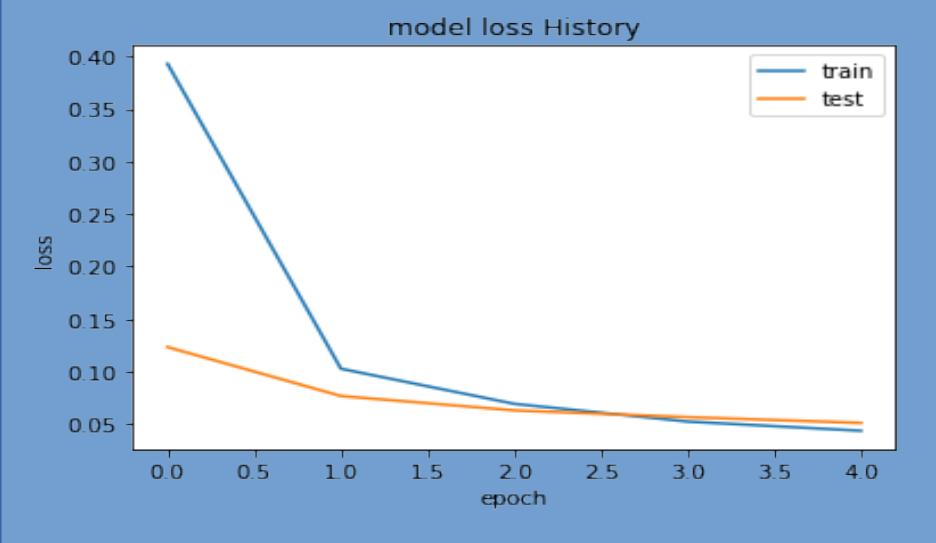
10000/10000 [=====] - 0s 30us/step
Test score: 8.15280571334064
Test accuracy: 97.51999974250793

Applied AI and IMAGE PROCESSING

CONV2:LENET MNIST HANDWRITTEN DIGIT RECOGNITION

PARAMETERS

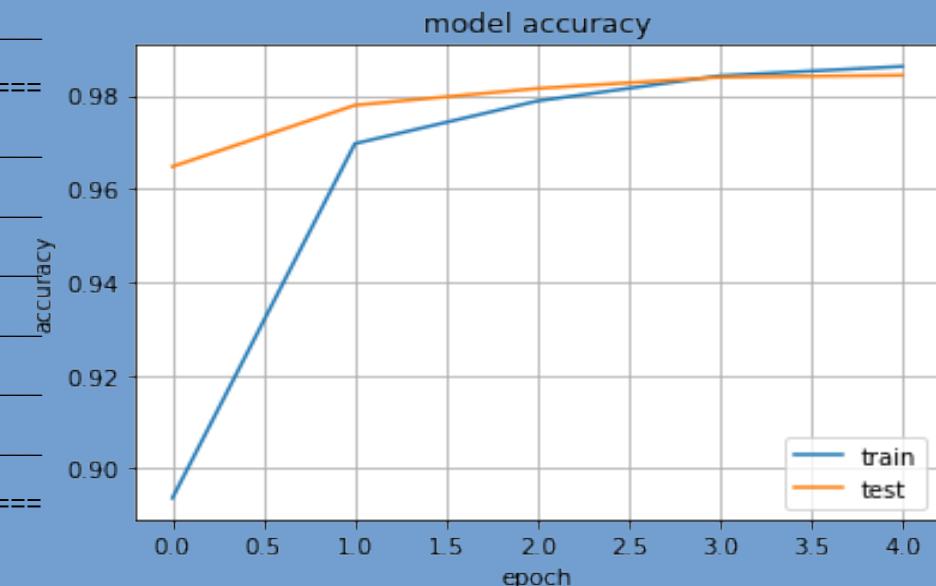
train samples 60000
test samples 10000
classes 10
hidden layers 128
dropout 0.25
learning rate 0.0002
batch_size 128
Validation split 0.2



Model: "sequential_1"

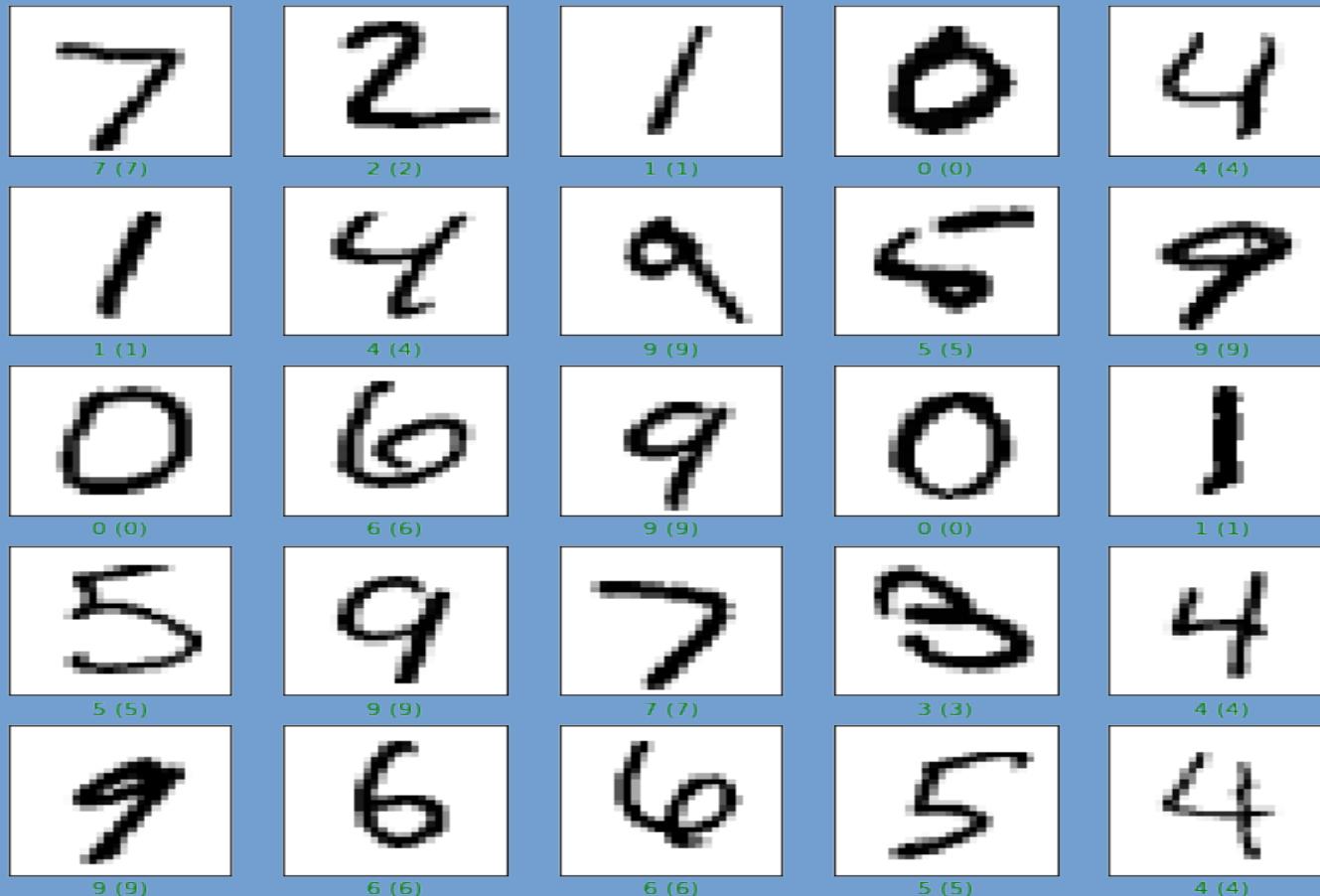
Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 20, 28, 28)	520
max_pooling2d_1 (MaxPooling2D)	(None, 20, 14, 14)	0
conv2d_2 (Conv2D)	(None, 50, 14, 14)	25050
max_pooling2d_2 (MaxPooling2D)	(None, 50, 7, 7)	0
flatten_1 (Flatten)	(None, 2450)	0
dense_1 (Dense)	(None, 500)	1225500
dense_2 (Dense)	(None, 10)	5010

Total params: 1,256,080
Trainable params: 1,256,080
Non-trainable params: 0



Applied AI and IMAGE PROCESSING

CONV2:LENET MNIST HANDWRITTEN DIGIT RECOGNITION



10000/10000 [=====] - 8s 766us/step

Test score: 3.788638093601912

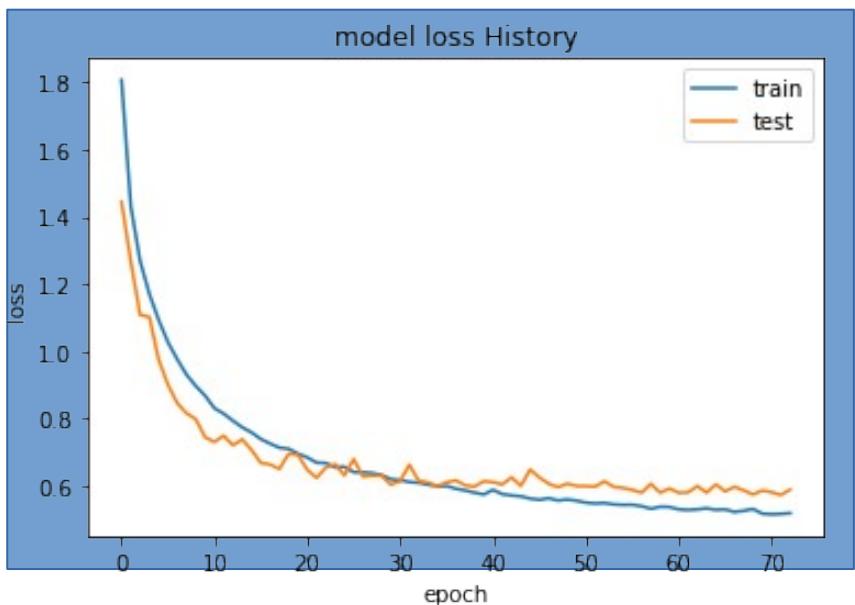
Test accuracy: 98.79000186920166

Applied AI and IMAGE PROCESSING

CONVNET : CIFAR10 CLASSIFICATION OF IMAGES

PARAMETERS

train samples 40000
test samples
10000
classes 10
learning rate 0.0005
batch_size 32
Validation split 0.2

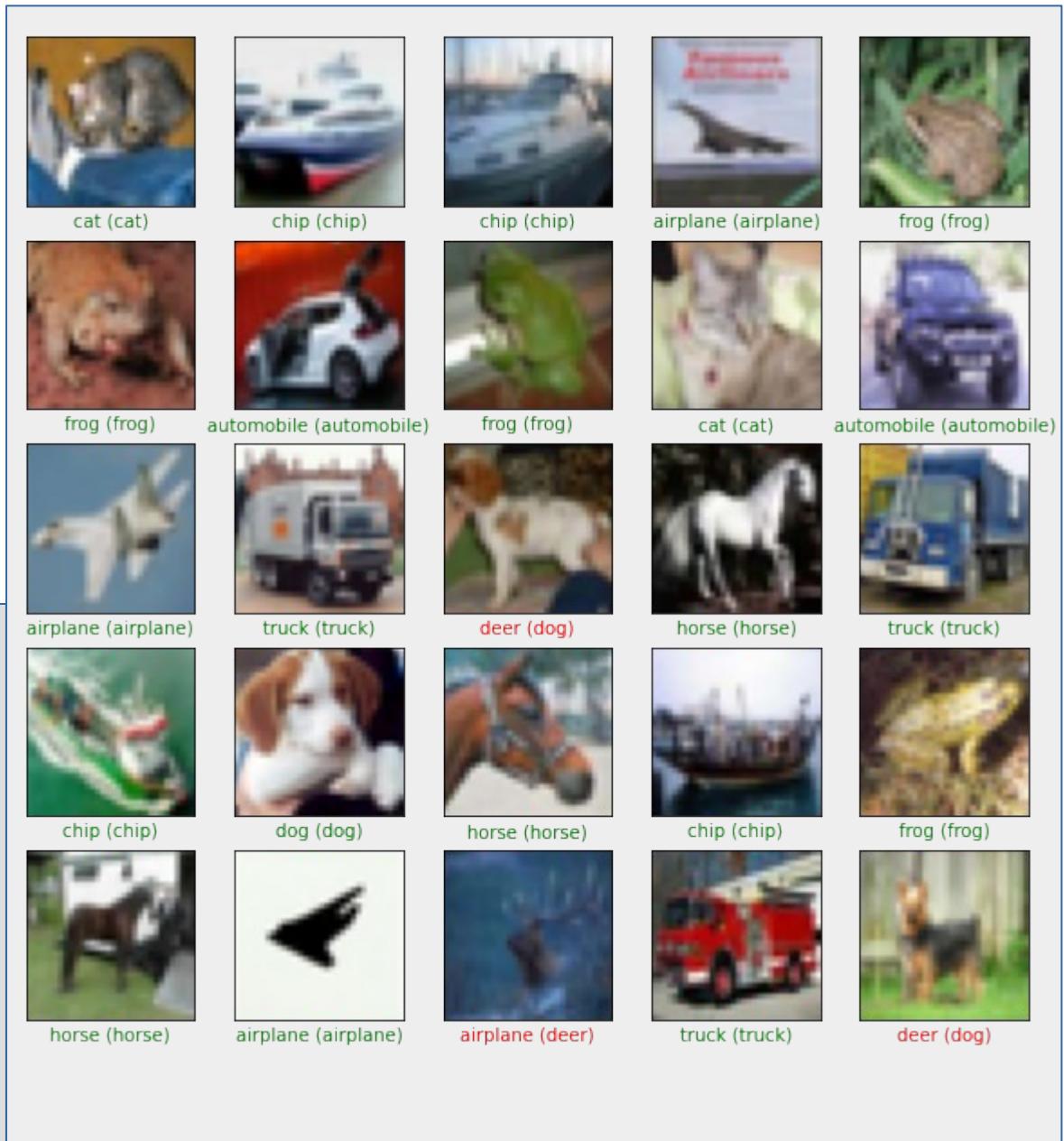
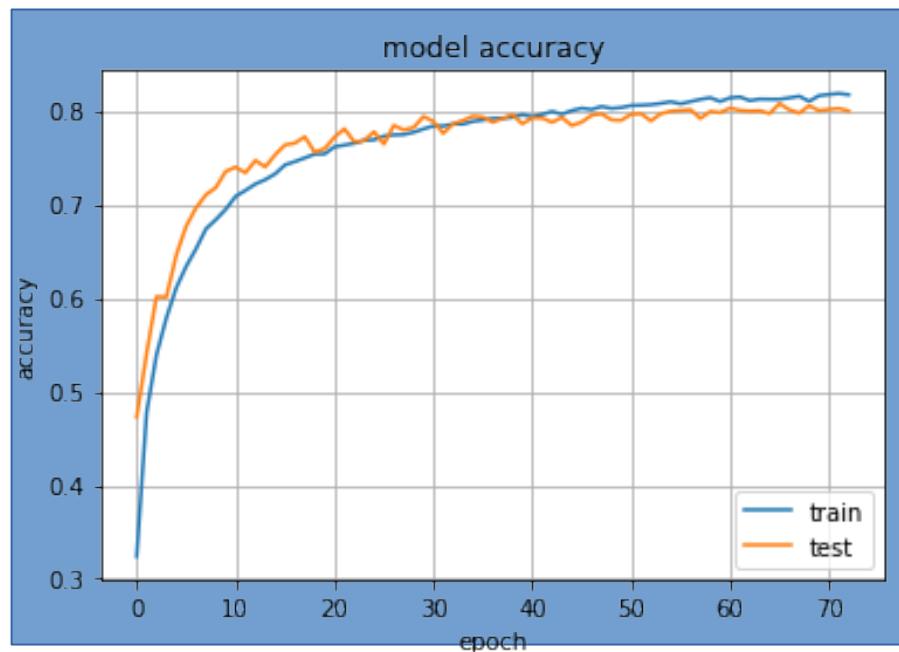


Model: "sequential_1"

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 32, 32, 16)	448
conv2d_2 (Conv2D)	(None, 32, 32, 16)	2320
max_pooling2d_1 (MaxPooling2D)	(None, 16, 16, 16)	0
dropout_1 (Dropout)	(None, 16, 16, 16)	0
conv2d_3 (Conv2D)	(None, 16, 16, 32)	4640
conv2d_4 (Conv2D)	(None, 14, 14, 32)	9248
max_pooling2d_2 (MaxPooling2D)	(None, 7, 7, 32)	0
dropout_2 (Dropout)	(None, 7, 7, 32)	0
conv2d_5 (Conv2D)	(None, 7, 7, 64)	18496
conv2d_6 (Conv2D)	(None, 5, 5, 64)	36928
max_pooling2d_3 (MaxPooling2D)	(None, 2, 2, 64)	0
dropout_3 (Dropout)	(None, 2, 2, 64)	0
flatten_1 (Flatten)	(None, 256)	0
dense_1 (Dense)	(None, 512)	131584
dropout_4 (Dropout)	(None, 512)	0
dense_2 (Dense)	(None, 10)	5130
activation_1 (Activation)	(None, 10)	0
<hr/>		
Total params: 208,794		
Trainable params: 208,794		
Non-trainable params: 0		

Applied AI and IMAGE PROCESSING

CONVNET : CIFAR10 CLASSIFICATION OF IMAGES



10000/10000

[==]-1s 87us/step

Test score:

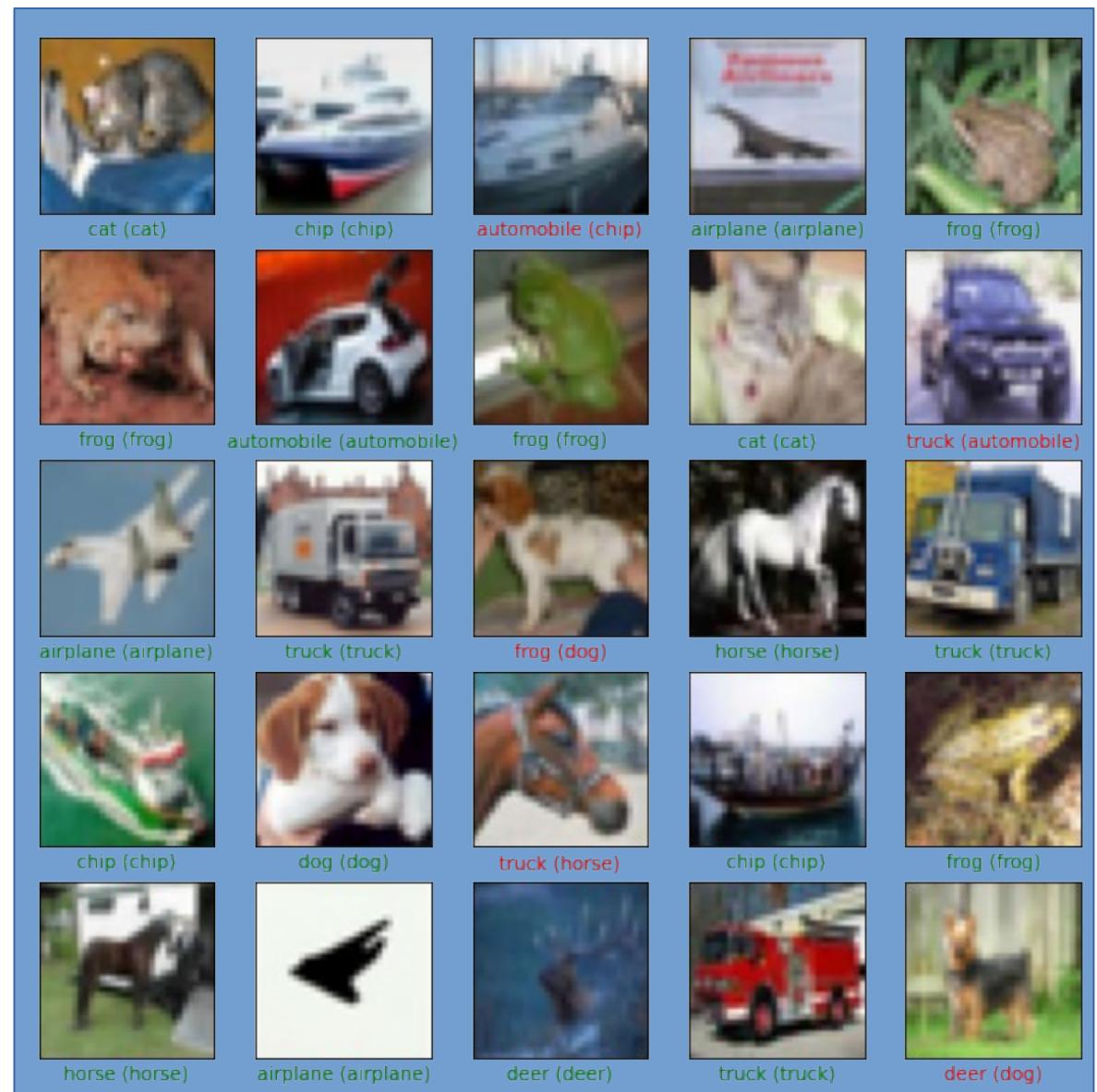
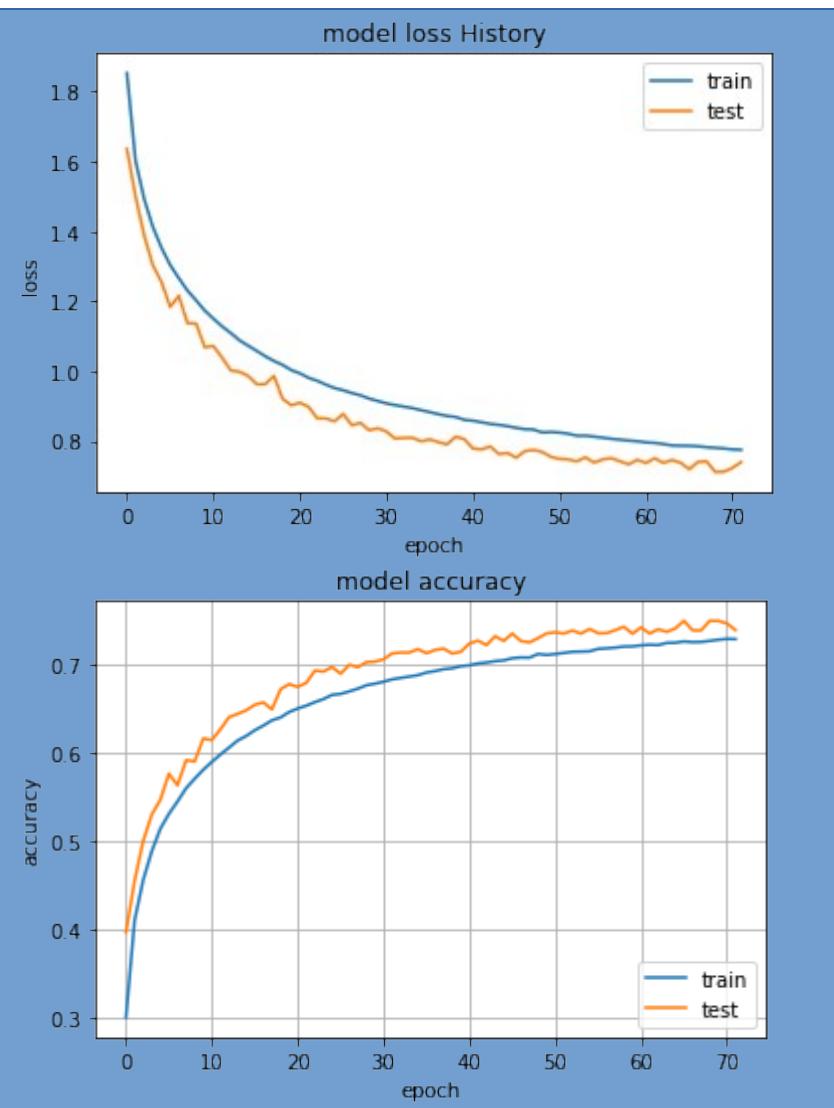
60.95441319942474

Test accuracy:

79.62999939918518

Applied AI and IMAGE PROCESSING

CONVNET : CIFAR10 CLASSIFICATION OF IMAGES
DATA AUGMENTEDx5 lr=2*10^-4 BATCH_SIZE=128



10000/10000 [=====] - 1s 113us/step

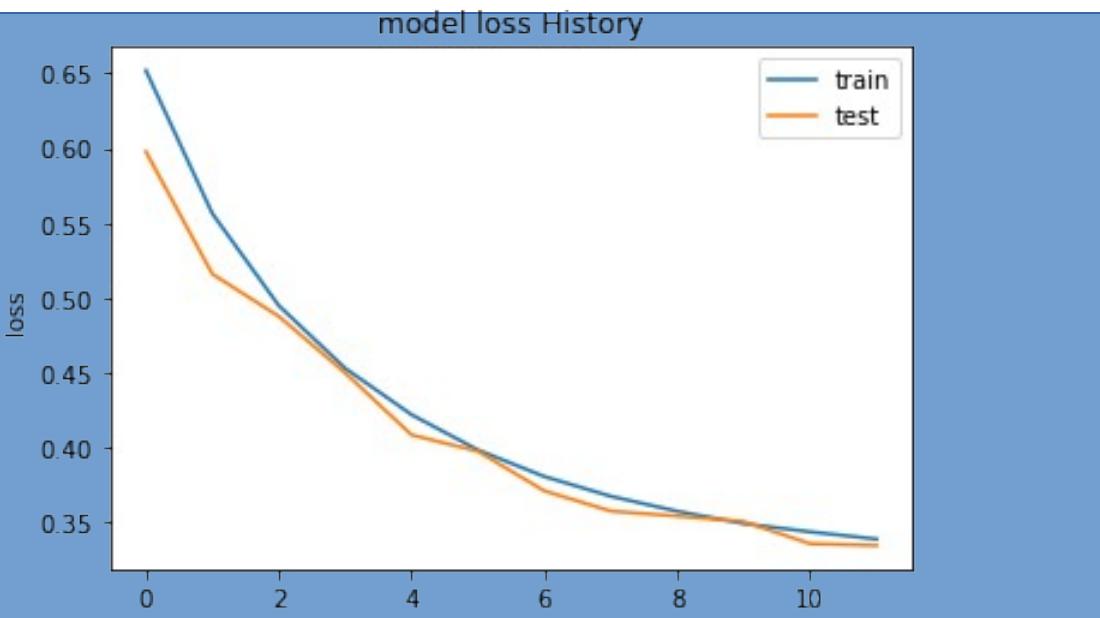
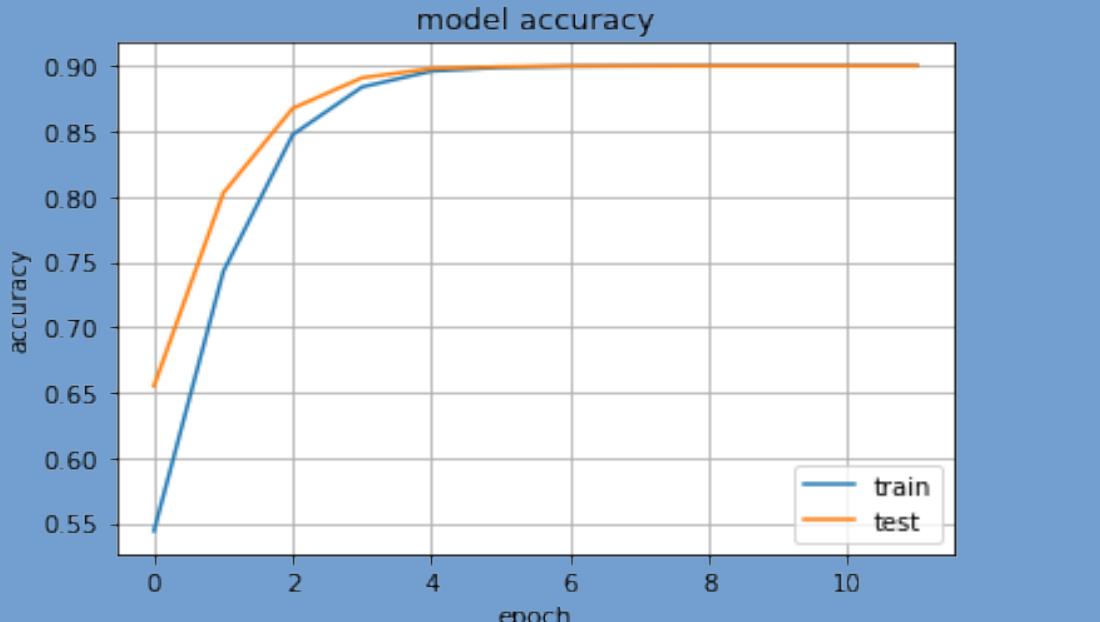
Test score: 82.35761503219604

Test accuracy: 71.90999984741211

Applied AI and IMAGE PROCESSING

CONVNET : STL10 CLASSIFICATION OF IMAGES

Image DataGenerator + PRE-TRAINED MODEL VGG16[-3 layers]



Model: "model_7"		
Layer (type)	Output Shape	Param #
input_4 (InputLayer)	(None, 96, 96, 3)	0
block1_conv1 (Conv2D)	(None, 96, 96, 64)	1792
block1_conv2 (Conv2D)	(None, 96, 96, 64)	36928
block1_pool (MaxPooling2D)	(None, 48, 48, 64)	0
block2_conv1 (Conv2D)	(None, 48, 48, 128)	73856
block2_conv2 (Conv2D)	(None, 48, 48, 128)	147584
block2_pool (MaxPooling2D)	(None, 24, 24, 128)	0
block3_conv1 (Conv2D)	(None, 24, 24, 256)	295168
block3_conv2 (Conv2D)	(None, 24, 24, 256)	590080
block3_conv3 (Conv2D)	(None, 24, 24, 256)	590080
block3_pool (MaxPooling2D)	(None, 12, 12, 256)	0
block4_conv1 (Conv2D)	(None, 12, 12, 512)	1180160
block4_conv2 (Conv2D)	(None, 12, 12, 512)	2359808
block4_conv3 (Conv2D)	(None, 12, 12, 512)	2359808
block4_pool (MaxPooling2D)	(None, 6, 6, 512)	0
block5_conv1 (Conv2D)	(None, 6, 6, 512)	2359808
block5_conv2 (Conv2D)	(None, 6, 6, 512)	2359808
block5_conv3 (Conv2D)	(None, 6, 6, 512)	2359808
block5_pool (MaxPooling2D)	(None, 3, 3, 512)	0
flatten_7 (Flatten)	(None, 4608)	0
dense_19 (Dense)	(None, 128)	589952
dense_20 (Dense)	(None, 128)	16512
dense_21 (Dense)	(None, 10)	1290
Total params: 15,322,442		
Trainable params: 5,327,370		
Non-trainable params: 9,995,072		

Applied AI and IMAGE PROCESSING

PRE-TRAINED MODELS PREDICTION :CLASSIFICATION OF IMAGES

FRACTALS IN NATURE

ResNet50

VGG16

VGG19

InceptionV3

aloe-spiral



'coil'

99.28587675094604

'coil'

85.86033582687378

'coil'

99.2819607257843

'coil'

84.0158104896545

broccoli-1



'broccoli'

66.43117070198059

'cauliflower'

43.12864542007446

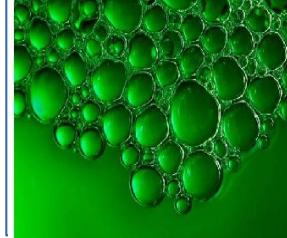
'cauliflower'

85.71616411209106

'cauliflower'

91.8545663356781

air-bubbles-2



'bubble'

45.6044465303421

'chain_mail'

31.28640353679657

'necklace'

31.665799021720886

'honeycomb'

75.99433660507202

pinecone2



'sea_cucumber'

27.89953649044037

'honeycomb'

25.767436623573303

'acorn'

84.96872782707214

'pineapple'

96.2532997131347

fractals



'chambered_nautilus'

99.94205236434937

'chambered_nautilus'

99.7738242149353

'chambered_nautilus'

99.53280687332153

'chambered_nautilus'

99.82097148895264

Applied AI and IMAGE PROCESSING

FACENET: FACE RECOGNITION TASK

FAmous MATHEMATICIANS

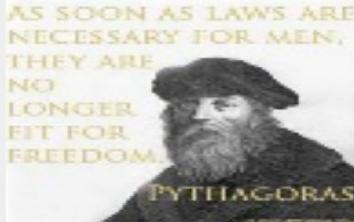
Bernhard-Riemann



Christiaan-Huygens



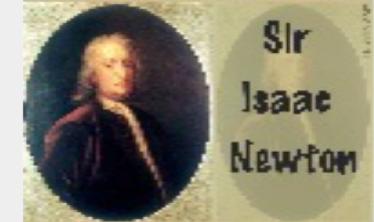
Pitagoras



Andrey-Kolmogorov



isaac-newton



Andrew-Wiles



CHristiaan-Huygens



Vladimir-Arnold



Stephen-Hawking



marjorie-lee-browne



Benoit_Mandelbrot



Leonardo-da-Vinci



euphemia-haynes



Srinivasa-Ramanujan

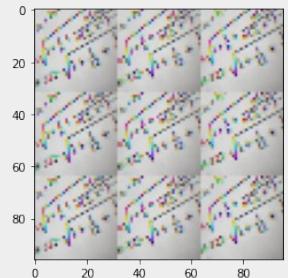


Gottfried-Leibniz

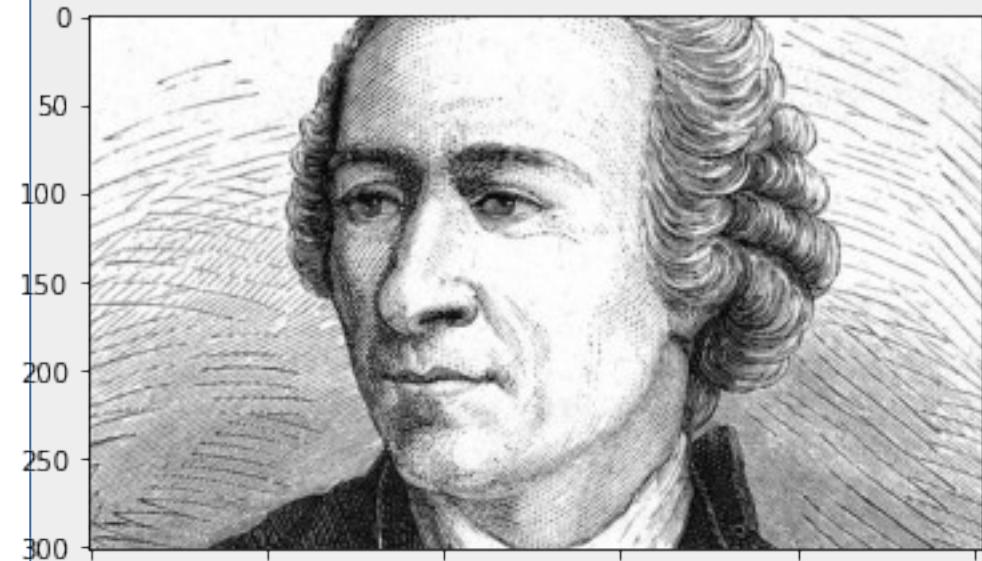


Applied AI and IMAGE PROCESSING

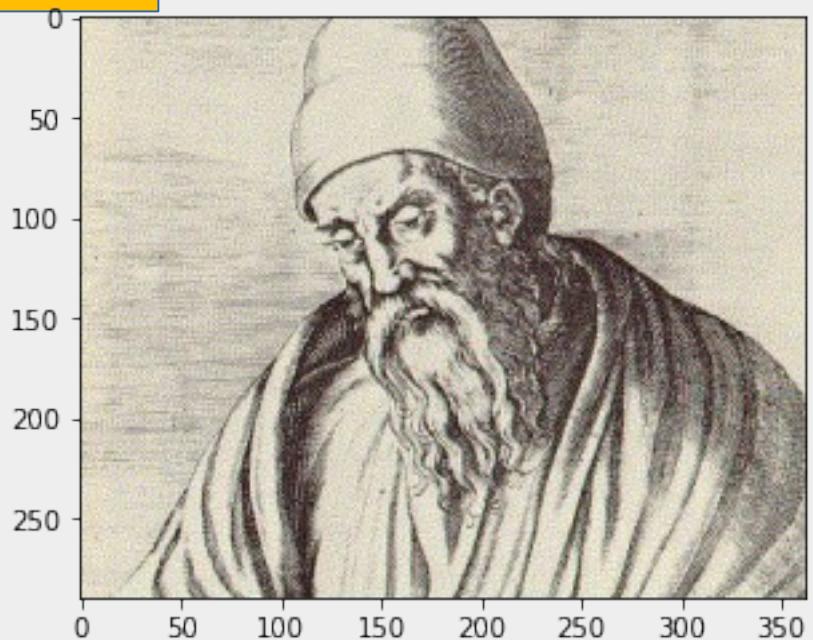
FACENET: FACE RECOGNITION TASK FAMOUS MATHEMATICIANS



It's not Rene-Descartes, please
have to continue studying
Out[28]:(0.9089534, False)



```
who_is_it("images96/Euler.jpg", database,  
FRmodel)  
It's Euler, the distance is 0.0
```



```
who_is_it("images96/Euclid.jpg", database,  
FRmodel)  
It's Euclid, the distance is 0.0
```

Applied AI and NLP

NATURAL LANGUAGE PROCESSING

Predict tags on StackOverflow with linear models

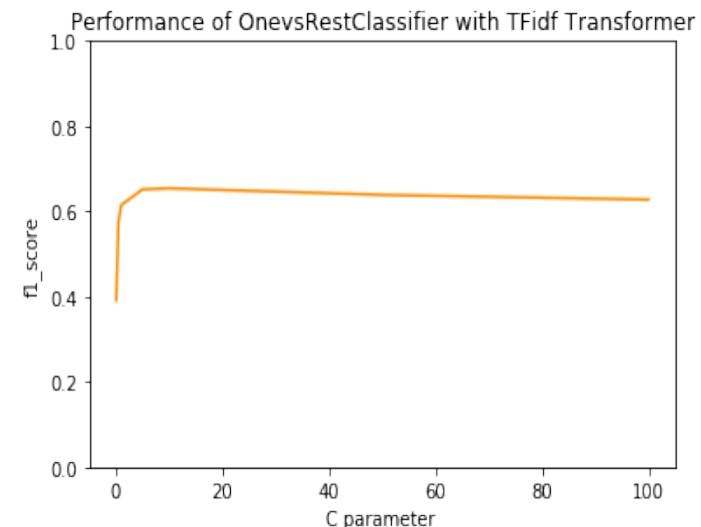
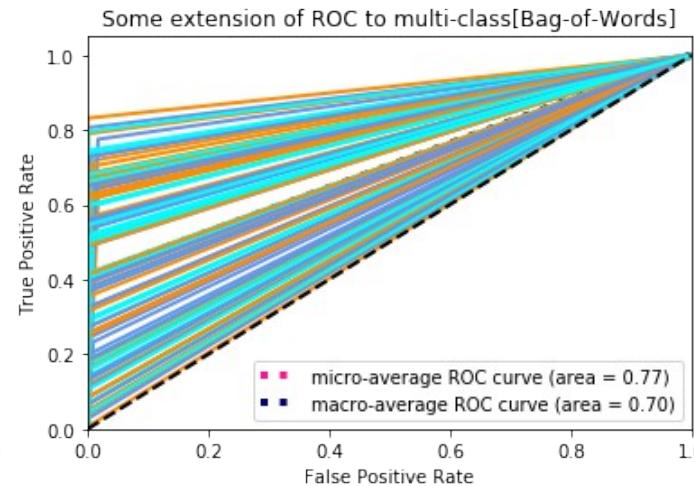
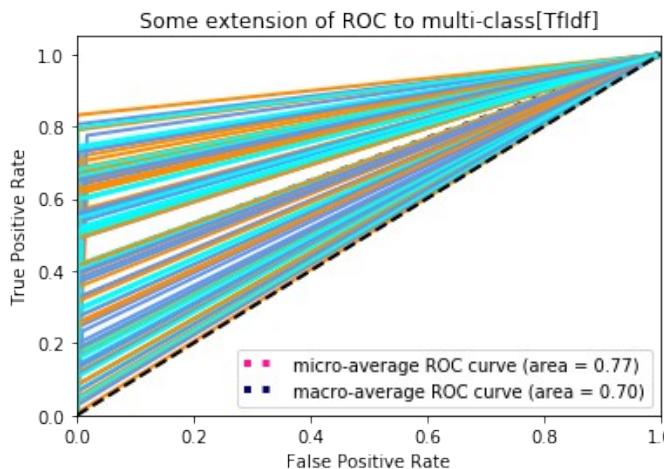
PARAMETERS

DICT_SIZE	5000
Validation split	0.2
Testing split	0.13
Bag of Words	5000
TF-IDF	18300
num_classes	100

LogisticRegression
 Bag-of-words
 0.6514944359375615
LogisticRegression
 TfIdf 0.6540109906751699
 Support Vector Machine
 Bag-of-words
 0.6581060304888663
 Support Vector Machine
 TfIdf 0.6614169003852886
MultinomialNB
 Bag-of-words
 0.6296148179299079
MultinomialNB
 TfIdf 0.398129431834211

Task Multilabel Classification:

0	mysql,php
1	javascript
2	python
3	javascript,jquery
4	android,java
5	php,xml
6	json
7	java,swing



Applied AI and NLP

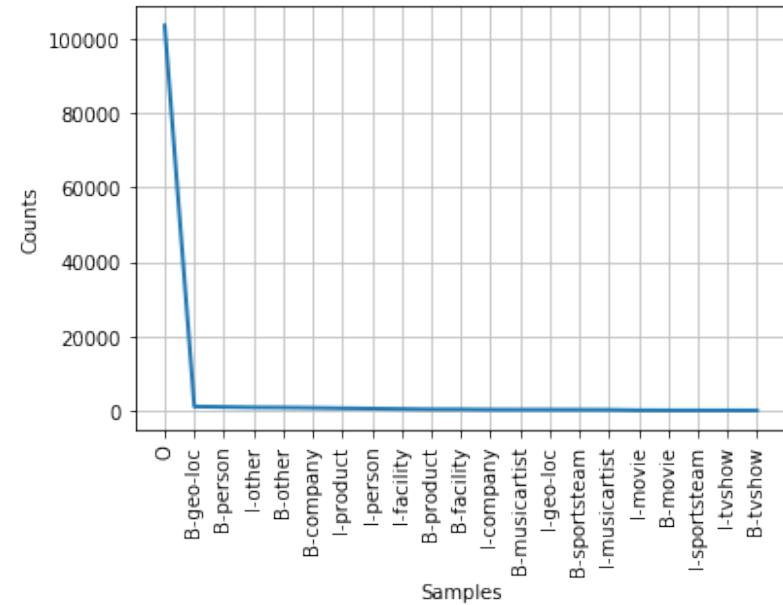
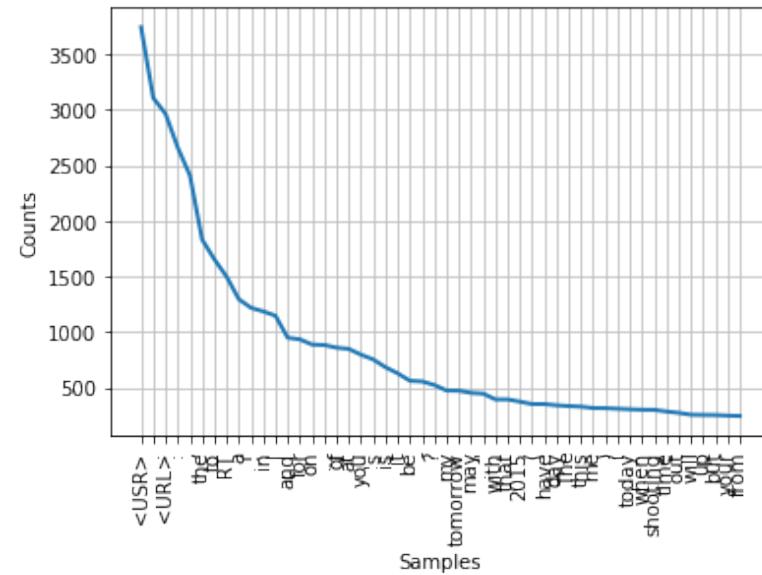
NATURAL LANGUAGE PROCESSING

Recognize named entities on Twitter with LSTMs

B-PER I-PER B-PRODUCT I-PRODUCT B-ORG I-ORG

PARAMETERS

vocabulary_size	112095
words_tokens	20503
sequence_length	41
embedding_dim	200
n_hidden_rnn	200
PAD_index	1
batch_size	32
n_epochs	4
learning_rate	0.0075
learning_rate_decay	2.0
dropout_keep_probability	0.7



----- Train set quality: -----

```

processed 105778 tokens with 4489 phrases; found: 4651 phrases; correct: 4174.
precision: 89.74%; recall: 92.98%; F1: 91.33
    company: precision: 91.72%; recall: 94.71%; F1: 93.19; predicted: 664
    facility: precision: 89.44%; recall: 91.72%; F1: 90.57; predicted: 322
    geo-loc: precision: 93.10%; recall: 97.59%; F1: 95.29; predicted: 1044
    movie: precision: 50.00%; recall: 57.35%; F1: 53.42; predicted: 78
    musicartist: precision: 80.57%; recall: 85.78%; F1: 83.09; predicted: 247
    other: precision: 88.96%; recall: 93.66%; F1: 91.25; predicted: 797
    person: precision: 93.54%; recall: 96.50%; F1: 95.00; predicted: 914
    product: precision: 89.33%; recall: 92.14%; F1: 90.71; predicted: 328
    sportsteam: precision: 85.07%; recall: 86.64%; F1: 85.84; predicted: 221
    tvshow: precision: 61.11%; recall: 37.93%; F1: 46.81; predicted: 36
  
```

Applied AI and NLP

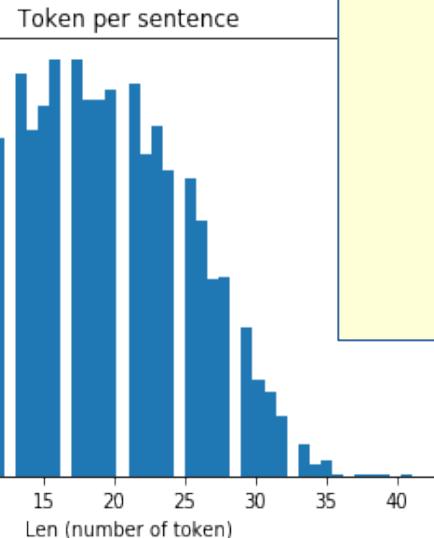
NATURAL LANGUAGE PROCESSING

Recognize named entities on Twitter with LSTMs

B-PER I-PER B-PRODUCT I-PRODUCT B-ORG I-ORG

PARAMETERS

vocabulary_size	112095
words_tokens	20503
sequence length	41
embedding_dim	200
n_hidden_rnn	200
PAD_index	1
batch_size	32
n_epochs	4
learning_rate	0.0075
learning_rate_decay	2.0
dropout_keep_probability	0.7



----- Validation set quality: -----
 processed 12836 tokens with 537 phrases; found: 401 phrases; correct: 192.
 precision: 47.88%; recall: 35.75%; F1: 40.94
 company: precision: 59.34%; recall: 51.92%; F1: 55.38; predicted: 91
 facility: precision: 46.15%; recall: 35.29%; F1: 40.00; predicted: 26
 geo-loc: precision: 64.84%; recall: 52.21%; F1: 57.84; predicted: 91
 movie: precision: 0.00%; recall: 0.00%; F1: 0.00; predicted: 3
 musicartist: precision: 25.00%; recall: 14.29%; F1: 18.18; predicted: 16
 other: precision: 40.62%; recall: 32.10%; F1: 35.86; predicted: 64
 person: precision: 45.31%; recall: 25.89%; F1: 32.95; predicted: 64
 product: precision: 13.79%; recall: 11.76%; F1: 12.70; predicted: 29
 sportsteam: precision: 23.53%; recall: 20.00%; F1: 21.62; predicted: 17
 tvshow: precision: 0.00%; recall: 0.00%; F1: 0.00; predicted: 0

----- Test set quality: -----
 processed 13258 tokens with 604 phrases; found: 526 phrases; correct: 254.
 precision: 48.29%; recall: 42.05%; F1: 44.96
 company: precision: 50.00%; recall: 46.43%; F1: 48.15; predicted: 78
 facility: precision: 51.35%; recall: 40.43%; F1: 45.24; predicted: 37
 geo-loc: precision: 73.44%; recall: 56.97%; F1: 64.16; predicted: 128
 movie: precision: 0.00%; recall: 0.00%; F1: 0.00; predicted: 5
 musicartist: precision: 9.52%; recall: 7.41%; F1: 8.33; predicted: 21
 other: precision: 39.00%; recall: 37.86%; F1: 38.42; predicted: 100
 person: precision: 55.81%; recall: 46.15%; F1: 50.53; predicted: 86
 product: precision: 10.53%; recall: 14.29%; F1: 12.12; predicted: 38
 sportsteam: precision: 27.27%; recall: 29.03%; F1: 28.12; predicted: 33
 tvshow: precision: 0.00%; recall: 0.00%; F1: 0.00; predicted: 0

Applied AI and NLP

NATURAL LANGUAGE PROCESSING

Learn to calculate with seq2seq model¶

Dataset of generated equations with the operators addition and substraction

PARAMETERS

```
vocab_size          15
embeddings_size     20
max_iter            7
hidden_size         512
start_symbol_id=word2id["^"]
end_symbol_id=word2id["$"]
padding_symbol_id=word2id["#"]
batch_size = 128
n_epochs = 10
learning_rate = 0.01
dropout_keep_probability = 0.5
max_len =20
```

ENCODER - DECODER ARCHITECTURE

Train: epoch 10

Epoch: [10/10],step: [1/625],loss: 1.120122

Epoch: [10/10], step: [201/625], loss: 1.166973

Epoch: [10/10], step: [401/625], loss: 1.115895

Epoch: [10/10], step: [601/625], loss: 1.134763

Test: epoch 10 loss: 1.1185263

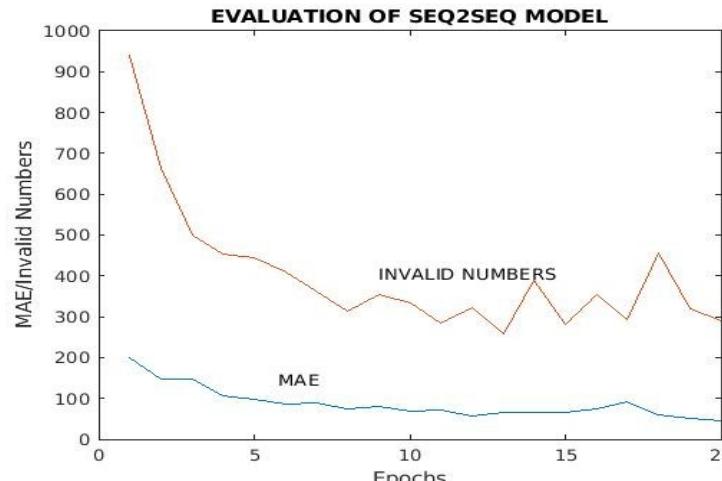
X: 8160-8424\$
Y: -264\$#
O: -288\$#

X: 1249+7575\$
Y: 8824\$#
O: 8887\$#

X: 565-567\$##
Y: -2\$###
O: -111\$#

EVALUATE RESULTS

```
Epoch: 1, MAE: 284.882726, Invalid numbers: 1232
Epoch: 2, MAE: 151.816467, Invalid numbers: 494
Epoch: 3, MAE: 122.000563, Invalid numbers: 454
Epoch: 4, MAE: 119.092943, Invalid numbers: 472
Epoch: 5, MAE: 112.519414, Invalid numbers: 555
Epoch: 6, MAE: 97.322286, Invalid numbers: 387
Epoch: 7, MAE: 102.530941, Invalid numbers: 382
Epoch: 8, MAE: 76.065245, Invalid numbers: 305
Epoch: 9, MAE: 75.231431, Invalid numbers: 344
Epoch: 10, MAE: 69.462707, Invalid numbers: 318
```



Applied AI and NLP

NATURAL LANGUAGE PROCESSING

TEXT-SUMMARIZATION WITH AMAZON REVIEWS WITH ATTENTION LAYER

MODEL (NO-FILTERING) PARAMETERS

54,018,847

vocab_size	56957
embeddings_size	500
max_len_text	80
max_len_summary	10
encoder	LSTM
decoder	LSTM
batch_size	64
rnn_len	256
n_layers	2
learning_rate	0.005
dropout_keep_probability	0.75
max_len =20	

MODEL PARAMETERS

FILTERING + EMB GLOV.6B.50d

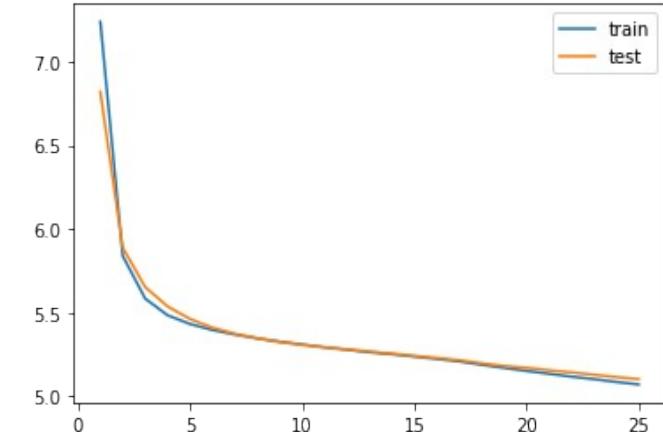
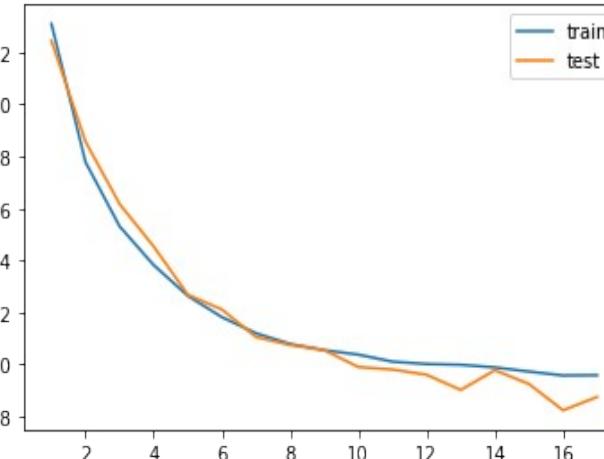
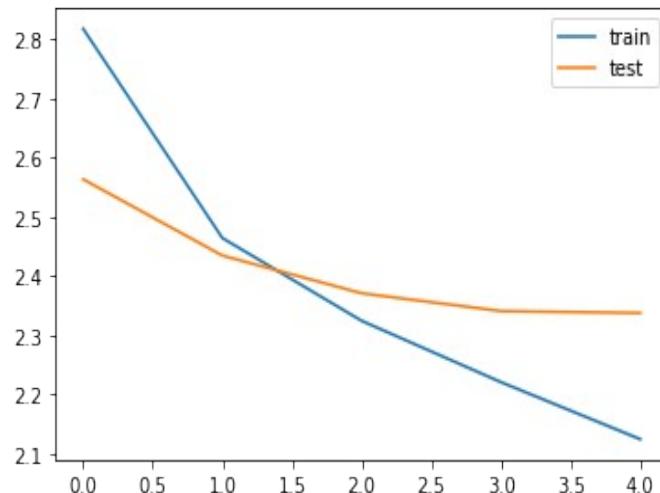
len_word_emb	36805
text/sum filtered	27125
embeddings_size	50
encoder/decoder	GRUCELL+ DROP + BDIRECT-DYNAMIC RNN
batch_size	64
rnn_len	128
n_layers	2
learning_rate	0.005
dropout_keep_probability	0.75

MODEL PARAMETERS

FILTERING + EMB CC.EN.300.VEC

len_word_emb	43127
text/sum filtered	27128
embeddings_size	300
encoder/decoder	GRUCELL+ DROP + BDIRECT-DYNAMIC RNN
batch_size	64
rnn_len	128
n_layers	2
learning_rate	0.00005
dropout_keep_probability	0.75

LOSS FUNCTIONS

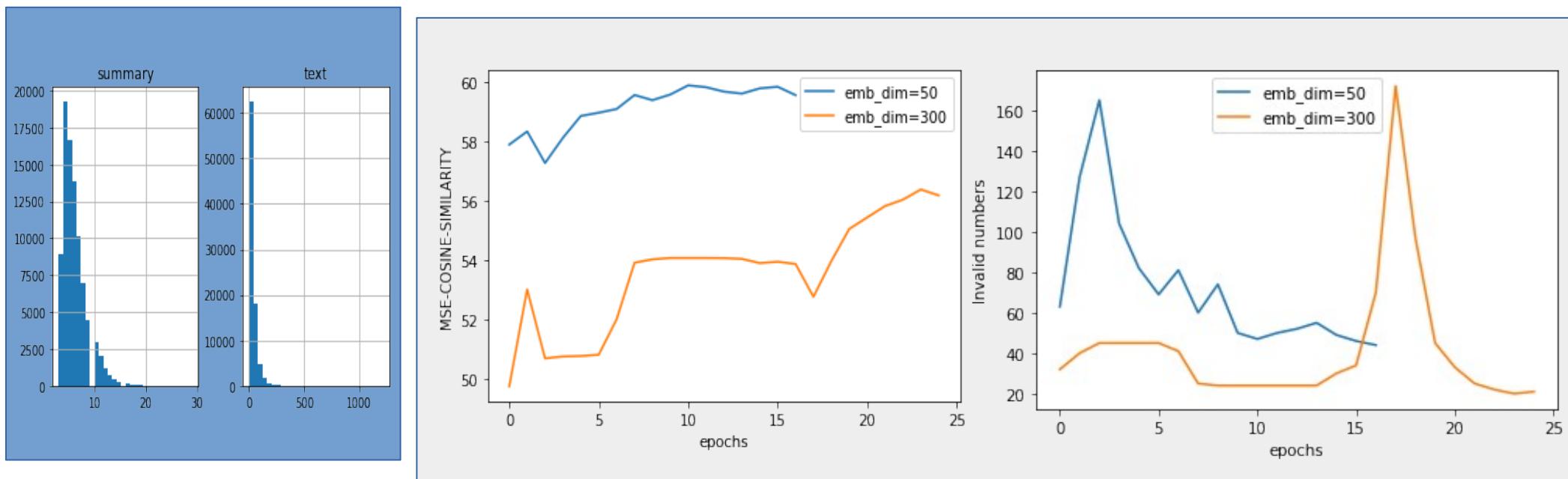


Applied AI and NLP

NATURAL LANGUAGE PROCESSING

TEXT-SUMMARIZATION WITH AMAZON REVIEWS

EVALUATION OF RESULTS



Embedding_dim=50

Text

Word Ids: [0, 721, 179, 6983, 399, 106, 460, 781, 2291, 3503, 5983, 1528, 5320, 359, 664, 71, 1738, 2802, 1646, 36803]
 Input Words: good strong coffee modified cup pack works well keurig brewer suppose less wasteful cannot beat price especially
 subscribe save

Summary

Word Ids: [0, 179, 137, 71]
 Response Words: good coffee at price
 Ground Truth: good coffee good price

Applied AI and NLP

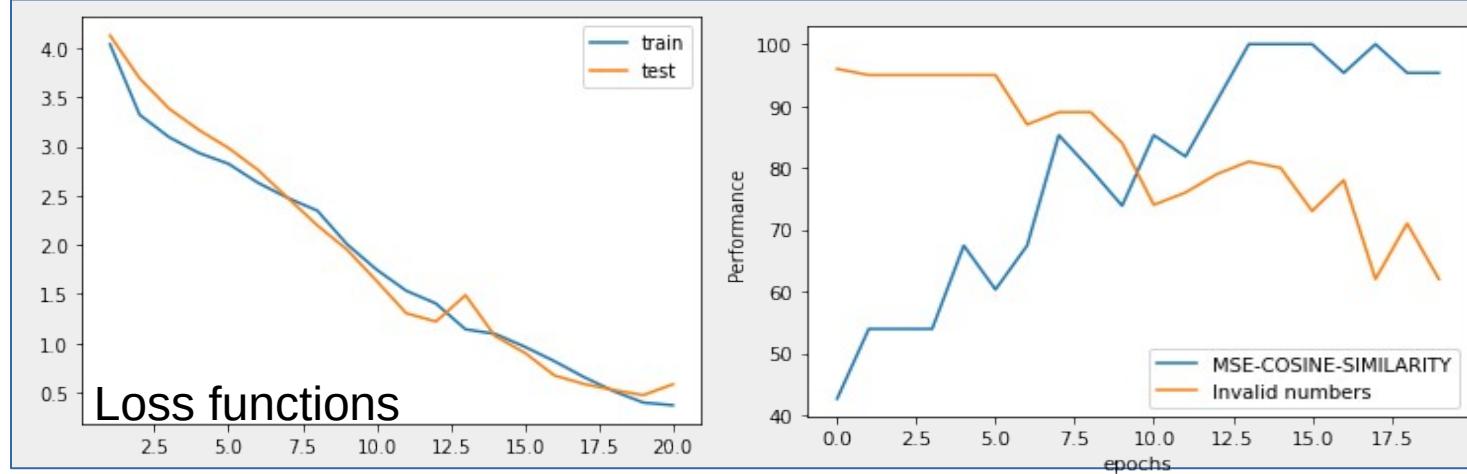
NATURAL LANGUAGE PROCESSING

ENGLISH-CATALÀ TRANSLATION

USING FILTER + PRE-TRAINED EMBEDDINGS

PARAMETERS

eng vocab_size	991
Cat vocab_size	1166
embeddings_size	300
FILTERED WORDS	78.71
rnn_len	256
n_layers	1
batch_size	32
learning_rate	0.0025
dropout_keep_probability =	0.75



Input Language;
index to word
mapping

11 ----> too
12 ----> late
7 ----> .
3 ----> <EOS>
3 ----> <EOS>

Target Language;
index to word
mapping

12 ----> massa
13 ----> tard
11 ----> .
3 ----> <EOS>

Text
Word Ids: [0, 305, 47, 81, 618, 213, 73, 619, 7, 3, 3]
Input Words: tell me your plans for the future .

Summary
Word Ids: [692, 56, 123, 693, 694, 149]
Response Words: explica m els teus plans per
Ground Truth: explica m els teus plans per al futur .

Text
Word Ids: [0, 16, 15, 70, 644, 189, 16, 645, 49, 7, 3, 3]
Input Words: i can t remember where i bought it .

Summary
Word Ids: [82, 724, 210, 19, 71, 306]
Response Words: no recorde on el vaig comprar
Ground Truth: no recordo on el vaig comprar .

Text
Word Ids: [0, 632, 67, 633, 80, 634, 635, 7, 3, 3]
Input Words: bangkok is thailand s capital city .

Summary
Word Ids: [712, 40, 98, 713]
Response Words: bangkok es la capital
Ground Truth: bangkok es la capital de tailandia .

Applied AI and NLP

NATURAL LANGUAGE PROCESSING

StackOverflow assistant bot

Dataset size 40000
Test size 11%

Intent recognizer
Tag classifier

Test accuracy = 0.9917954545454546
Test accuracy = 0.800625

RANKING QUESTIONS WITH EMBEDDINGS

'c#':	394451
'java':	383456
'javascript':	375867
'php':	321752
'c_cpp':	281300
'python':	208607
'r':	36359
'ruby':	99930
'vb':	35044
'swift':	34809

