Introduction to Statistics on HPC Using Matlab

THE UNIVERSITY • OF ARIZONA



- Another Introduction to HPC workshop
- Leverages Matlab Statistics Toolbox
- Examples shown on HPC







Statistics and AI

Statistics plays a crucial role in AI by providing the foundation for data analysis, modeling, and decision-making.

Within AI, statistics serves as the backbone for various tasks such as predictive modeling, pattern recognition, and data interpretation.

By leveraging statistical techniques, AI systems can analyze complex data sets, identify patterns, and make informed decisions based on empirical evidence.





David Blackwell

Nvidia has named their new GPU after David Blackwell, an American mathematician and statistician.

Blackwell provided pioneering work in Bayesian statistics. Bayesian statistics is a powerful framework that uses probability theory to quantify uncertainty and update beliefs based on new evidence.

In AI applications, Bayesian methods play a crucial role in tasks such as machine learning, pattern recognition, and probabilistic reasoning.









THE UNIVERSITY

OF ARIZONA

A

The Language of Technical Computing

- Campus license
- HPC license

Why Use Matlab on HPC?

- 1. Use for large datasets
- 2. Use for multiple concurrent jobs
- 3. Collaboration workspaces
- 4. It looks very similar to laptop



Statistics and Machine Learning Toolbox

- Use descriptive statistics, visualizations and clustering.
- Fit probably distributions to data.
- Generate random numbers for:
 - Monte Carlo simulations
 - Regression and Classification to draw inferences
 - Principal Component Analysis, Regularization,
 - Dimensionality Reduction and Feature Selection



ARIZONA HPC Systems Apps - Files - Jobs - Clusters - Interactive Apps - 🗇 My Interactive Sessions

This is the UArizona *Open OnDemand* server **Please NOTE:** "windfall" jobs will be restarted or terminated without notice if pre-empted by a "standard" job in queue.

OnDemand provides an integrated, single access point for all of your HPC resources.

Pinned Apps A featured subset of all available apps



Interactive HPC:

ood.hpc.arizona.edu





HE UNIVERSITY

OF ARIZONA

Matlab GUI

This app will launch the Matlab GUI on a UAz cluster. You will be able to interact with the Matlab GUI through a Turbo VNC based session.

Cluster

ElGato Cluster

□ I would like to receive an email when the session starts

Matlab GUI version

Default

This defines the version of the application you want to use

Run Time



Enter the number of wall clock hours the job is allowed to run.

Core count on a single node



Enter the number of cores on a single node that the job is allowed to use.



\$



Session was successfully created.

Home / My Interactive Sessions

Interactive Apps	Matlab GUI (1819802)		1 node 2 cores Running
Desktops			
☐ Interactive Desktop	Host: >_cpu39.elgato.hpc.arizona.edu		😢 Delete
Uls	Created at: 2024-03-04 17:08:39 MST		
🛦 Abaqus GUI	Time Remaining: 3 hours and 59 minute	S	
Ansys Workbench GUI	Session ID: a6a7c86d-0790-4507-97ff-	780f234d4251	
Mathematica GUI	Compression	Image Quality	
🛦 Matlab GUI	0 (low) to 9 (high)	0 (low) to 9 (high)	
	Launch Matlab GUI		View Only (Share-able Link)
n Stata GUI			
VSCode GUI			











Statistics and Machine Learning Toolbox Analyze and model data using statistics and machine learning

- Descriptive Statistics
- Cluster Analysis ANOVA
- Regression and Classification
- Dimensionality Reduction
- Probability Distributions
- Hypothesis Tests
- Industrial Statistics
- Analysis of Big Data
- Code Generation





Efficiently summarize and validate large datasets using descriptive statistics and frequency distributions, including measures of central tendency and dispersion, stem plots, bar/pie charts and histograms.

Of the available examples we choose Visualizing Multivariate Data using three different statistical plots.





	i 🛃 🗼		MATLAB R2023a - academic use						C	_ • ×	
	НОМЕ	PLOTS /	APPS		2 10	🖥 🤋 🖉 🔁 😨) 💽 Sear	ch Docume	ntation 🔎	🌲 Christop	oher S 👻
	New New Script Live Script	New Open	Find Files Compare	Import Data	Clean Data VA	o Variable ▼ Save Workspace Clear Workspace RIABLE	e V CODE	SIMULINK		RESOURCES	Ā
Cut and	* * 🖬 🗖 🕅	🖨 / 🕨 home 🕨 i	u13 🕨 chri:	sreidy 🕨							+ p
Julanu	Current Folder	\odot	Comman	d Window	N		-0				\odot
ste:	E Name 4		New to M.	ATLAB? Se	e resour	ces for <u>Getting Sta</u>	arted.				X
Open Arrow Select	□ □ □ Using □ □ bayes Using □ □ chapel >> □ □ cuda >> □ □ □ □			j 2 thread(s) on compute node. j 2 thread(s) on compute node. pad carbig					Evaluate Selec Open Selectior Help on Select	F9 F4 F1	
Clipboard	data		<pre>X = [MPG,Acceleration,Displacement,Weight,Horsepony varNames = {'MPG'; 'Acceleration'; 'Displacement';</pre>					Function Brows	Shift+F1		
Paste with		lipboard	14-22-1						Function Hints		Ctrl+F1
Right click	git X = [MR bello-warNam	rbig PG,Acceleration,Displ les = {'MPG'; 'Acceler	acement,Wei ration'; 'Displ	ight,Horsepo lacement'; '\	ower]; Weight'; 'H	orsepower'};			Cut Copy Paste		Ctrl+W Alt+W Ctrl+Y
Paste	* ht / this fall	der is not on your -click hello-worl	MATLAB p d to make i	ath. It your cur	rent fold	er or	1	-	Select All Find		Ctrl+X, H Ctrl+H
	Details	"Add to Path" from	n its contex (Do	kt menu to not show (b add it t this mes	age again)	Clear	J	Print Print Selection Page Setup	m	264-2551501094
	Workspace	۲						-	Clear Comman	d Window	Ctrl+L
THE UNIVERSITY OF ARIZONA	Name L Acceleration	Value 406×1 douk •									

📣 MATLAB



- 1. Open Arroy
- Select 2. Clipboard
- 3. Paste with cmd-V
- **Right click** 4.
- Paste 5.













OF ARIZONA







Cluster Analysis K-Means Clustering



Data often naturally fall into groups (or clusters) of observations, where the characteristics of objects in the same cluster are similar.

This K-means example is extensive, and we just demonstrate the first features.





Cluster Analysis K-Means Clustering







Cluster Analysis K-Means Clustering

HE UNIVERSITY

OF ARIZONA









In an ANOVA model, each grouping variable represents a fixed factor. The levels of that factor are a fixed set of values. The goal is to determine whether different factor levels lead to significantly different response values or outcomes.

This example starts the ANOVA process. It continues with F-Statistics and Variance Components using mean squares and confidence bounds





ANOVA Analysis of Variance

ommand Wind	ow			O
New to MATLAB? S	See resource	s for <u>Getting Sta</u>	ted.	×
>> load mile >> factory >> carmod = >> mileage =	eage = repmat(1 [ones(3,3) = mileage(:	:3,6,1); ; 2*ones(3,3));	1;	
correctory = ta	<pre>actory(:); amod(:);</pre>			
[mileage fac	tory carmo	d]		
	(275)			
ans =				
33.3000	1.0000	1.0000		
33,4000	1.0000	1.0000		
32.9000	1.0000	1.0000		
32,6000	1.0000	2.0000		
32,5000	1.0000	2,0000		
33.0000	1.0000	2.0000		
34.5000	2.0000	1.0000		
34.8000	2,0000	1.0000		
33.8000	2.0000	1.0000		
33,4000	2.0000	2.0000		
33.7000	2,0000	2,0000		
33.9000	2.0000	2.0000		-
37.4000	3.0000	1.0000		
36.8000	3.0000	1.0000		
37.6000	3.0000	1.0000		
1			2000	



Part

1



ANOVA Analysis of Variance

ans =						-			
33, 3000	1.0000	1.0000							
33,4000	1.0000	1.0000							
32.9000	1.0000	1.0000							
32.6000	1.0000	2.0000		Figure	1- N-Way	ANOVA	1.000		1
32.5000	1.0000	2.0000	<u> </u>						
33.0000	1.0000	2.0000	<u>File Edit V</u> iew <u>I</u> nsert	<u>T</u> ools <u>D</u>	esktop	<u>W</u> indow <u>H</u> el	0		ಿ
34.5000	2.0000	1.0000		Analys	sis of '	Variance			
34.8000	2.0000	1.0000	Source	- Sum So	d f	Mean Sa	E	Droha	
33.8000	2.0000	1.0000				nean oq.			T
33.4000	2.0000	2.0000	Factory	53.3511	2	26,6756	1333.78	0.0007	
33.7000	2.0000	2.0000	Car Model	1.445	ī	1.445	72.25	0.0136	
33.9000	2.0000	2.0000	Factory:Car Model	0.04	2	0.02	0.18	0.8411	
37.4000	3.0000	1.0000	Error	1.3667	12	0.1139			
36.8000	3.0000	1.0000	Total	56.2028	17				
37.6000	3.0000	1.0000							
36.6000	3.0000	2.0000		nstrained (Type III)	sums of square	20		1.0
37.0000	3.0000	2.0000		Shotramed (type my	sams or squar			-
36.7000	3.0000	2.0000				300			

A



Statistics and Machine Learning Toolbox Regression Bayesian Analysis for a Logistic Regression Model

Use the Regression Learner app or programmatically train and assess models such as linear regression, Gaussian processes, support vector machines, neural networks, and ensembles.

This example starts on Bayesian inferences for a logistic regression model. Try the whole example on your own.





Statistics and Machine Learning Toolbox Regression **Bayesian Analysis for a Logistic Regression Model**



OF ARIZONA

Statistics and Machine Learning Toolbox Regression **Bayesian Analysis for a Logistic Regression Model**



OF ARIZONA





Statistics and Machine Learning Toolbox Classification Train Classification Ensemble

Use the Classification Learner app or programmatically train and validate models such as logistic regression, support vector machines, boosted trees, and shallow neural networks.

This example shows how to create a classification tree ensemble.





Statistics and Machine Learning Toolbox Classification Train Classification Ensemble



IE UNIVERSITY

OF ARIZONA



Statistics and Machine Learning Toolbox Classification Train Classification Ensemble







Statistics and Machine Learning Toolbox Dimensionality Reduction and Feature Extraction

This toolbox has many tools including:

- Sequential feature selection
- Partial least squares regression
- Principal components regression
- Principal components analysis
- Discriminant Analysis Classifier
- Random forest predictors
- Feature extraction workflow
- Extract mixed signals
- Visualize high-dimensional data
- Factor analysis

UNIVERSITY ARIZONA

- Nonnegative matrix factorization
- Multidimensional scaling



Statistics and Machine Learning Toolbox Dimensionality Reduction and Feature Extraction Classical Multidimension Scaling

Command Window

New to MATLAB? See resources for <u>Getting Started</u>.
>> rng default; % For reproducibility

```
X = [normrnd(0,1,10,3),normrnd(0,.1,10,1)];
D = pdist(X,'euclidean');
>> [Y,eigvals] = cmdscale(D);
>> format short g
[eigvals eigvals/max(abs(eigvals))]
```

ans =

35.41	1
11.158	0.31511
1.6894	0.04771
0.1436	0.0040553
5.4637e-15	1.543e-16
3.2157e-15	9.0813e-17
2.2212e-15	6.2727e-17
1.3179e-15	3.7219e-17
-2.3377e-15	-6.6019e-17
-3.47e-15	-9.7995e-17



Part



 \odot

×

Part 2

Statistics and Machine Learning Toolbox Dimensionality Reduction and Feature Extraction Classical Multidimension Scaling

```
>> maxerr4 = max(abs(D - pdist(Y))) % Exact reconstruction
  maxerr4 =
     1.7764e-15
  >> maxerr3 = max(abs(D - pdist(Y(:,1:3)))) % Good reconstruction in 3D
  maxerr3 =
       0.043142
  >>
  >> maxerr2 = max(abs(D - pdist(Y(:,1:2)))) % Poor reconstruction in 2D
  maxerr2 =
        0.98315
  >> max(max(D))
  ans =
         5.8974
fx >>
```









Fit continuous and discrete distributions, use statistical plots to evaluate goodness-of-fit, and compute probability density functions.

This example shows how to fit multiple probability distribution objects to the same set of sample data.





```
Command Window
```

New to MATLAB? See resources for Getting Started.

```
>> load carsmall
>> Origin = categorical(cellstr(Origin));
MPG2 = MPG(Origin~='Italy');
Origin2 = Origin(Origin~='Italy');
Origin2 = removecats(Origin2,'Italy');
>> [WeiByOrig,Country] = fitdist(MPG2,'weibull','by',Origin2);
[NormByOrig,Country] = fitdist(MPG2, 'normal', 'by', Origin2);
[LogByOrig,Country] = fitdist(MPG2,'logistic','by',Origin2);
[KerByOrig,Country] = fitdist(MPG2, 'kernel', 'by', Origin2);
>> WeiByOrig
WeiByOrig =
  1×5 cell array
  Columns 1 through 2
    {1×1 prob.WeibullDistribution}
                                      {l×1 prob.WeibullDistribution}
  Columns 3 through 4
                                      {1×1 prob.WeibullDistribution}
    {1×1 prob.WeibullDistribution}
  Column 5
    {1×1 prob.WeibullDistribution}
```



Part



 \odot

x



Part 2

OF ARIZONA



MATLAB

Part 3

HE UNIVERSITY

OF ARIZONA

Command Window

{'France' }

{'Germany'} {'Japan' }

{'Sweden']

>> WeiUSA = WeiByOrig{5};

NormUSA = NormByOrig{5};

LogUSA = LogByOrig{5};

{'USA'

New to MATLAB? See resources for Getting Started.

KerUSA = KerByOriq{5}; x = 0:1:50;pdf Wei = pdf(WeiUSA,x); pdf Norm = pdf(NormUSA,x); pdf Log = pdf(LogUSA,x); pdf Ker = pdf(KerUSA,x); >> data = MPG(Origin2=='USA'); figure histogram(data,10,'Normalization','pdf','FaceColor',[1, >> line(x,pdf Wei, 'LineStyle', '-', 'Color', 'r') line(x,pdf Norm,'LineStyle','-.','Color','b') line(x,pdf_Log,'LineStyle','--','Color','g') line(x,pdf Ker,'LineStyle',':','Color','k') legend('Data','Weibull','Normal','Logistic','Kernel','L title('MPG for Cars from USA') xlabel('MPG') fx :==



Part 4





Draw inferences about a population based on statistical evidence from a sample. Perform parametric and nonparametric tests.

This example shows how to use hypothesis testing to analyze gas prices measured across the state of Massachusetts during two separate months.





Current Folder	💿 Command Window		
🗋 Name 🛆	New to MATLAB? See resources for Getting Star	ed.	
8 🗀pycache 9 🗀 bayes 9 🗀 chapel 8 🗀 cuda 8 🗀 culinux 1 data	Using 2 thread(s) on compute node. Using 2 thread(s) on compute node.		
Bocuments dvml Clipboard evra git hello-world	j e2];	Evaluate Selection Open Selection Help on Selection	F9 F4 F1
inte	Cle	Function Browser Show Function Browse	Shift+F1 r Button
intre-to-hps		Function Hints	Ctrl+F1
inux		Cut	Ctrl+W
ietails		Сору	ALE+W
Varkonaca		Paste	Ctrl+Y
lame / Value		Select All	Ctrl+X, H
ans 2		Find	Ctrl+H
Ready		Print Print Selection	





OF ARIZONA





Command Window

New to MATLAB? See resources for Getting Started.

Using 2 thread(s) on compute node. Using 2 thread(s) on compute node. >> load gas prices = [pricel price2]; >> normplot(prices) Warning: MATLAB has disabled some advanced graphics rendering features by switching to software OpenGL. For more information, click <u>here</u>. >> lillietest(price1)

ans =

```
0
```

```
>> lillietest(price2)
```

```
ans =
```

```
Θ
```

```
>> sample_means = mean(prices)
```

```
sample_means =
```

115.1500 118.5000

fx >>



Part

2





 \odot



Part 3





MATLAB

Part 4

OF ARIZONA



Industrial Statistics

Analyzing Survival or Reliability Data Survivor Functions for Two Groups Hazard and Survivor Functions for Different Groups Cox Proportional Hazards Model Object Cox Proportional Hazards for Censored Data Cox Proportional Hazards with Time-Dependent Covariates **Control Charts**







Industrial Statistics Control Charts

Control charts are created with the <u>controlchart</u> function. Any of the following chart types may be specified:

- •Xbar or mean
- Standard deviation
- Range
- •Exponentially weighted moving average
- Individual observation
- Moving range of individual observations
- Moving average of individual observations
- Proportion defective
- Number of defectives
- •Defects per unit
- Count of defects





Industrial Statistics Control Charts











