# **SOCR 2022 MDP Project Summaries**

The one-page summaries below describe the main SOCR MDP R&D Projects for 2022 (January-December) https://www.socr.umich.edu/html/SOCR\_Research.html

### GDrive: <u>https://drive.google.com/drive/folders/1qvFrrak3hn8irPSoG1675vVIDBF6IIOr</u> GSlides: <u>https://docs.google.com/presentation/d/1r1GTx6KgFwLOuu0cDsGxppwN48T\_e0WYtfpqOj7wu6E/</u> SOCR Project Leaders:

• Programming:

Simeone Marino, Alex Kalinin, Ivo Dinov

- Methods (CBDA, GrayRain/VH, DataSifter): Simeone Marino
- Analytics:
- Spacekime Analytics:

Simeone Marino, Ivo Dinov Yueyang Shen, Ivo Dinov

### SOCR Trainees/Students:

...TBD...

## **Project Summaries**

| Project Area                    | Skills                                | Likely Majors                       |
|---------------------------------|---------------------------------------|-------------------------------------|
| Programming Subteam: SOCRAT     | UI/UX design, HTML5, JavaScript,      | Computer Science (CSE/CS-LSA)       |
| (Charts, Wrangler, Modeler,     | Adobe Illustrator, Canvas             | School of Information (SI)          |
| Analyses, Tools) (2-3 students) |                                       |                                     |
| TensorFlow.JS                   | https://js.tensorflow.org             | Computer Science (CSE/CS-LSA)       |
|                                 | https://js.tensorflow.org/api/latest/ |                                     |
| UKBB t-SNE, BrainViewer         | https://codepen.io/pen?&editors=1011  |                                     |
| Methods (CBDA & DataSifter)     | Technical math background, R-         | Math, CS, Eng, Physics, Stats, STEM |
| (4 students)                    | computing                             |                                     |
| DataSifter & CBDA & GrayRain/VH |                                       |                                     |
| AI/ML Methods                   | Develop Reinforcement Learning,       | Math, CS, Eng, Physics, Stats, STEM |
| (2-4 students)                  | and deep networks methods             |                                     |
| Analytics                       | R/Python, statistical modeling, high- | Statistics, Biostatistics,          |
| (4 students)                    | throughput data analytics, machine    | Bioinformatics                      |
| <u>TDA</u>                      | learning                              | Math                                |
| Biomed/Health Applications (see |                                       | Computer Science (CSE/CS-LSA)       |
| <u>Case-Studies</u> )           |                                       |                                     |
| Spacekime Analytics (sub-team – | Information measures, entropy KL      | Physics, math or engineering        |
| working directly with the PI)   | divergence, PDEs, Dirac's bra-ket     | background is preferred             |
| (4 students)                    | operators.                            |                                     |
| www.spacekime.org               | See <u>The Enigmatic Kime: Time</u>   |                                     |
|                                 | Complexity in Data Science at the     |                                     |
|                                 | University of Michigan Institute for  |                                     |
|                                 | Data Science (MIDAS) Seminar Series,  |                                     |
|                                 | Slidedeck, YouTube video of this      |                                     |
|                                 | seminar                               |                                     |

### SOCR Computing servers:

- ARC-TS: <u>https://arc-ts.umich.edu/open-ondemand/</u>
- SOCR-pipeline: socr-pipeline.nursing.umich.edu
- SOCR-RShiny: rshiny.umms.med.umich.edu
- SOCR-Lighthouse: <u>https://lighthouse.arc-ts.umich.edu</u> (Lighthouse User Guide)

# SOCR 2022 MDP Project: SOCRAT

### SOCR Project Leaders: Alex Kalinin / Ivo Dinov

| Website:          | https://socr.umich.edu/HTML5/SOCRAT/                                     |
|-------------------|--|
| GitHub:           | https://github.com/SOCR/SOCRAT   |
| Training Modules: | https://github.com/SOCR/socr-tutorials                                   |
| GDrive:           | https://drive.google.com/drive/folders/1UrNpNDI5sWoXW61YwP02NSv3PBbxfvpC |

### Description

The Statistics Online Computational Resource Analytics Toolbox (SOCRAT) is a Dynamic Web Toolbox for Interactive Data Processing, Analysis, and Visualization. It's purely built using HTML5 standards and JavaScript (core library) as well as node.js,

### **Student Skills**

- EECS, Computer Science (CSE/CS-LSA) and School of Information (SI)
- UI/UX design, HTML5, JavaScript

### **Project Goals**

- Go through the Training Modules, practice HTML/JS/Angular/Node programming
- Get your GitHub domain going and pull current SOCRAT branch
- Choose 1-2 deliverables, go over current design, start expansion, include unit tests, pilot development
- Coordinate with team

### Deliverables

- Expanded collection of Charts
- Expanded collection of Data-Modelers
- Expanded collection of (parametric and non-parametric) Statistical Analyses
- Expanded collection of machine learning classification, prediction, clustering and analytics modules.

### **Team Activities**

- Weekly team Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

### References

- Review the websites
- Alexandr A. Kalinin, Selvam Palanimalai, and Ivo D. Dinov. 2017. SOCRAT Platform Design: A Web Architecture for Interactive Visual Analytics Applications. In Proceedings of HILDA'17, Chicago, IL, USA, May 14, 2017, 6 pages. <u>DOI:10.1145/3077257.3077262</u>

**IDE** for development (Eclipse, WebStorm, IntelliJ, Netbeans, ..., RStuio, Spyder/Py)

# SOCR 2022 MDP Project: Methods: CBDA

### SOCR Project Leaders: Simeone Marino

| Website:          | http://socr.umich.edu/HTML5/CBDA/  |
|-------------------|--|
| GitHub:           | https://github.com/SOCR/CBDA   |
| C-RAN Package:    | https://cran.r-project.org/web/packages/CBDA                             |
| Training Modules: | http://socr.umich.edu/HTML5/CBDA/  |
| GDrive:           | https://drive.google.com/drive/folders/1hjwtgz64A_IUsnRK1gv7mGSJ3HdBHaRW |

### Description

The SOCR Compressive Big Data Analytics (CBDA) Project conducts research and implements efficient computational algorithms to tackle the Big Data problems of representation and analysis of complex heterogeneous information. Big Data cannot be loaded and processed as a whole. CBDA implements a real-time efficient divide-and-conquer strategy to deconstruct the Big Data into meaningful pieces of information that can be eventually reconstructed for actionable knowledge and predictive analytics.

### **Student Skills**

- Probability, stats, math, numerical methods, optimization
- R programming with RStudio (IDE) experience

### **Project Goals**

- Go through the provided materials and references
- Download the CBDA Package
- Practice with test-cases (<u>https://umich.instructure.com/courses/38100/files/folder/Case\_Studies</u>)
- Identify specific R&D direction to go deeper into an meaningfully contribute to CBDA
- Coordinate with team

### Deliverables

- New CBDA methods
- Expanded collection of machine learning forecasting, prediction, classification, clustering methods to expand the available CBDA algorithms
- Release new versions of CBDA R package and publish CBDA #2 manuscript
- Python/Perl scripts to speed up the subsampling strategy with Big Data > 100Gb-1Tb

### **Team Activities**

- Weekly team face-to-face/Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

- Review the websites
- Marino S, Xu J, Zhao Y, Zhou N, Zhou Y, Dinov, ID. (2018) Controlled feature selection and compressive big data analytics: Applications to biomedical and health studies, PLoS ONE 13(8): e0202674, DOI: 10.1371/journal.pone.0202674.
- Marino, S, Zhao, Y, Zhou, N, Zhou, Y, Toga, AW, Zhao, L, Jian, Y, Yang, Y, Chen, Y, Wu, Q, Wild, J, Cummings, B, Dinov, ID. (2020). Compressive Big Data Analytics: An ensemble meta-algorithm for highdimensional multisource datasets, PLoS ONE, 15(8):e0228520, DOI: 10.1371/journal.pone.0228520.

# SOCR 2022 MDP Project: Methods: DataSifter

### SOCR Project Leaders: Simeone Marino

| Website:          | http://DataSifter.org  |
|-------------------|--|
| GitHub:           | https://github.com/SOCR/DataSifter                                       |
| C-RAN Package:    | (lite version pending)   |
| Training Modules: | http://DataSifter.org  |
| GDrive:           | https://drive.google.com/drive/folders/1jVT5pTa_n8xHjUszn1u5gwTzyvPLtszj |

### Description

The SOCR DataSifter is a novel method, and an efficient R package, for on-the-fly de-identification of structured Clinical/Epic/PHI data. This approach provides complete administrative control over the risk of data identification when sharing large clinical cohort-based medical data. At the extremes, the data-governor may specify that either null data or completely identifiable data is generated and shared with the data-requester. This decision may be based on data-governor determined criteria about access level, research needs, etc. For instance, to stimulate innovative pilot studies, the data office may dial up the level of protection (which may naturally devalue the information content in the data), whereas for more established and trusted investigators, the data governors may provide a more egalitarian dataset that balances preservation of information content and sensitive-information protection.

### **Student Skills**

- Probability, stats, math, numerical methods, optimization
- R programming with RStudio (IDE) experience

### **Project Goals**

- Go through the provided materials and references
- Download the DataSifter-lite Package
- Practice with test-cases (<u>https://umich.instructure.com/courses/38100/files/folder/Case\_Studies</u>)
- Identify specific R&D direction to go deeper into and meaningfully contribute to DataSifter methods, implementation and/or validation
- Coordinate with team
- Examine the open-source Generative Pre-trained Transformer 3 (GPT-3), autoregressive language DL model, which has shown some impressive results in synthetic generation of human-like text/speech.
  - Try to integrate with DataSifterText (perhaps for the next paper)?
  - <u>https://github.com/openai/gpt-3</u>
  - <u>https://en.wikipedia.org/wiki/GPT-3</u>
  - https://arxiv.org/abs/2005.14165

### Deliverables

- Test this Python <u>SDV</u> package <u>https://ealizadeh.com/blog/sdv-library-for-modeling-datasets</u> <u>https://sdv.dev/SDV/user\_guides/evaluation/index.html</u>
- Remember we use the R-package charlatan: <u>https://cran.r-project.org/web/packages/charlatan/index.html</u>
- See: <u>https://bmcmedresmethodol.biomedcentral.com/articles/10.1186/s12874-020-00977-1</u>
- Implement these new privacy metrics: https://sdv.dev/SDV/user\_guides/evaluation/single\_table\_metrics.html#privacy-metrics
- New DataSifter methods/algorithms (e.g., addressing text, time-varying, graph data organizations)
- Release new versions of DataSifter R package
- Coordinate/support collaborators

### **Team Activities**

- Weekly team face-to-face/Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

### References

• Review the websites

- Marino, S, Zhou, N, Zhao, Yi, Wang, L, Wu Q., and Dinov, ID. (2018) DataSifter: Statistical Obfuscation of Electronic Health Records and Other Sensitive Datasets, Journal of Statistical Computation and Simulation, pp: 1-23, DOI: 10.1080/00949655.2018.1545228.
- Zhou, N, Wang, L, Marino, S, Zhao, Y, Dinov, ID. (2022), <u>DataSifter II: Partially Synthetic Data Sharing of Sensitive Information Containing Time-varying Correlated Observations</u>, <u>Journal of Algorithms & Computational Technology</u>, 15:1–17, DOI: <u>10.1177/17483026211065379</u>.

# SOCR 2022 MDP Project: Webapps & Data Analytics

SOCR Project Leaders: Simeone Marino, Ivo Dinov

 Website:
 <many, e.g., <a href="http://socr.umich.edu/HTML5">http://socr.umich.edu/HTML5">http://socr.umich.edu/HTML5</a>

 GitHub:
 <a href="https://github.com/SOCR">https://github.com/SOCR/ALS\_PA</a>

 Training Modules:
 <a href="http://DSPA.predictive.space">http://DSPA.predictive.space</a>

 GDrive:
 <a href="https://drive.google.com/drive/folders/1sN1fLYA0oLf114e1REJRthaMD0jXBs7w">https://drive.google.com/drive/folders/1sN1fLYA0oLf114e1REJRthaMD0jXBs7w</a>

### Description

The SOCR Webapp & Data Analytics projects are focused on interrogating massive amounts of complex biomedical and health data. Each project tackles multiple case-studies using R/RMD/RStudio, RShiny Services, and Python/Jupyter Notebook and the SOCR-Flux Compute Server

(<u>https://docs.google.com/document/d/1UmBq\_BMiMeUcijvKUCzPeG3tKZaWkinVtKrVWenPK1Y</u>). The webapp development will use R markdown notebook, RShiny web-servers, and Google BigQuery Datasets.

### **Student Skills**

- Biostats, quantitative analytics, probability, stats, math, numerical methods, optimization
- R programming with RStudio (IDE) experience, and/or Python/Jupyter Notebook

### **Project Goals**

- Go through the provided materials and references
- Review the SOCR Data Analytics Publications (<u>http://socr.umich.edu/people/dinov/publications.html</u>)
- Review the SOCR R-environment (<u>https://drive.google.com/file/d/1-u9adsMIYmMkcPD9W\_6BbfC1IMETsHF\_/</u>)
- Practice with test-cases (<u>https://umich.instructure.com/courses/38100/files/folder/Case\_Studies</u>)
- Identify specific case-study and an R&D direction to go deeper into an meaningfully contribute
- Coordinate with team

### Deliverables

- New SOCR end-to-end data analytics protocols
- Analytical results, abstracts, publications, presentations, research findings, etc.
- MIMIC-III analytics
- Baby-growth and mother-obesity relations
- Data Value Metric (DVM)
- European Economics Indicators (longitudinal analytics)
- 2D, 3D, 4D Visualization of complex data
- Coordinate/support collaborators

### **Team Activities**

- Weekly team face-to-face/Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

- Review the websites and listed resources
- https://shiny.med.umich.edu/apps/dinov/RShinyApp\_PIPM/
- https://socr.shinyapps.io/RShinyApp\_PIPM/

# SOCR 2022 MDP Project: Data Analytics - MIMIC-III

SOCR Project Leaders: Simeone Marino, Ivo Dinov

## Website: TBD

GitHub: https://github.com/SOCR

### Training Modules:

- Data Science & Predictive Analytics: <u>http://DSPA.predictive.space</u>
- Previous SOCR Data Analytics Publications: <u>http://socr.umich.edu/people/dinov/publications.html</u>
- Gaining access to the dataset requires an online training module; see onboarding materials below <u>https://drive.google.com/drive/u/1/folders/1Y6Yqq1CuTkHQ5rZg-C9r8\_je18nM886l</u>

GDrive: https://drive.google.com/drive/folders/1sN1fLYA0oLf1I4e1REJRthaMD0jXBs7w

### Description

This SOCR Data Analytics project is focused on interrogating the MIMIC-III database, a large collection of ~43,000 critical care patients from an ICU in Boston, MA. We will use R/RStudio, Python/Jupyter, and the SOCR-Flux Compute Server<sup>1</sup> to digest the vital signs, laboratory results, free-text data, and waveforms available in this unique dataset and predict clinical outcomes via statistical modeling tools.

### <sup>1</sup>SOCR-Flux Compute server:

https://docs.google.com/document/d/1UmBq\_BMiMeUcijvKUCzPeG3tKZaWkinVtKrVWenPK1Y

### **Student Skills**

- Biostats, quantitative analytics, probability, stats, math, numerical methods
- Programming experience in R (with RStudio) or Python (with Jupyter Notebook)
- Relational databases & structured query language (SQL)

### **Project Goals**

- Review the provided materials and references (see above)
- Request access to the MIMIC-III dataset (<u>https://mimic.physionet.org/gettingstarted/access/</u>)
   This involves an online but comprehensive human subjects research ethics course
- Practice with demo dataset (<u>https://physionet.org/works/MIMICIIIClinicalDatabaseDemo/</u>) and the MIMIC Query Builder (<u>https://querybuilder-lcp.mit.edu/dashboard.cgi</u>)
- Identify specific research aims and questions of interest to the team
- Coordinate with team to create a reproducible, accessible answer to these specific aims

### Deliverables

- New SOCR end-to-end data analytics protocols
- Data extraction & time-alignment tools for the MIMIC-III dataset
- Build statistical models to predict meaningful clinical outcomes
- Analytical results, abstracts, publications, presentations, research findings, etc.
- Visualization of complex, multidimensional data

### **Team Activities**

- Weekly team face-to-face/Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

# SOCR 2022 MDP Project: SOCR TensorFlow/TensorBoard Apps

SOCR Project Leader: Peijing Li, Lyn(Lingzi) Liao, Alex Kalinin, Ivo Dinov

 Website:
 https://socr.umich.edu/HTML5/SOCR\_TensorBoard\_UKBB

 GitHub:
 https://github.com/SOCR/97-tensorflowjs-quick-start

 Training Modules:
 https://js.tensorflow.org/tutorials/

 GDrive:
 https://drive.google.com/drive/folders/1wJY8539tpLmYiJc\_vKZvI6oDVDAHTQu9

### Description

The SOCR TensofFlowJS/TensorBoardJS project aims to design, built, validate and release new webapps based on the ML TensorFlow framework. For example, students will dive deep into TensorFlow.JS (<u>https://js.tensorflow.org</u>, <u>https://js.tensorflow.org/api/latest/</u>, <u>https://codepen.io/pen?&editors=1011</u>) and TensorBoard.JS (<u>https://github.com/tensorflow/tensorboard</u>, <u>https://www.tensorflow.org/guide/summaries\_and\_tensorboard</u>).

### Student Skills

- EECS, Computer Science (CSE/CS-LSA) and School of Information (SI)
- AngularJS, TensorFlowJS, TensorBoard, JavaScript, HTML5

### **Project Goals**

- Go through the Training Modules, practice HTML/JS/Angular/Node programming
- Get your GitHub domain going and start pilot testing various applications
- Use SOCR Data to experiment
- Review Vegi's SOCR t-SNE TensorFlow Webapp (<u>http://socr.umich.edu/HTML5/SOCR\_TensorBoard\_UKBB</u>)
- Coordinate with team
- Rapid RDD (research, development and deployment) is needed in this project

### Deliverables

- 2-5 new SOCR TF/TB Apps
- ...

### **Team Activities**

- Weekly team Zoom meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

### References

• Review the websites

# SOCR 2022 MDP Project: Interactive Graphical Probability Distribution Calculator

### SOCR Project Leader: Ivo Dinov (Mark Bobrovnikov, Jared Chai)

 Website:
 http://Distributome.org\_&

 Mitps://Distributome.org\_&
 https://shiny.med.umich.edu/apps/dinov/SOCR\_DistribCalc\_RShiny\_App/

 GitHub:
 https://github.com/distributome

 Training Modules:
 https://github.com/SOCR/socr-tutorials\_& http://dspa.predictive.space/

 GDrive:
 https://drive.google.com/drive/folders/184p8VNSOumYEG\_SOxIo4MyLVtanq9xLY

 Cools Stand along
 DMD southaut) that address the peak of the shallenge

**Goal**: Stand-alone RMD source and a demo HTML (RMD-output) that address the goal of the challenge, provide the desired functionality, and implement it efficiently, without any back-end support (e.g., no shiny server apps). **Deliverables**: Stand-alone RMD source and a demo HTML (RMD-output) that address the goal of the challenge, provide the desired functionality, and implement it efficiently.

. Background:

- Review DSPA <u>Chapter 2 (for probability distributions)</u> and <u>Chapter 5 (for plot\_ly demos)</u>
- See this DSPA Interactive Normal Probability CDF calculation tool.
- Review, experiment and play with the <u>Probability Distributome Calculators</u>. Try at least 2-dozen distributions what works well and what can be improved there? Mind the selection of parameters and the choice for function to plot (PDF, CDF).

**Desired Functionality**: The schematic below illustrates the core functionality of the interactive probability distribution calculator. Be creative in your solution.



- Include a drop-down list for the user to select the distribution
- Include an effective strategy to specify the parameters of the selected distribution, mind that the parameter number, interpretation and values will be different for different distributions.
- Make sure you keep the interactive aspect of the interface (*plot\_ly* style interactivity)
- Make sure all axes are appropriately scaled, labeled and drawn.

### **Functionality Annotations**

- A: x and y axes ranges and labels
- B: Appropriate title (placed to avoid overlaps
- C: Selection of the specific distribution should include at least 20 distributions, see class notes the Distributome Calculators.
- D: appropriate section of the specific distribution parameters

https://www.socr.umich.edu/html/SOCR\_Research.html

- E: cut off of the critical value (Z)
- **F**: Ranges of the animation slider should match the x-axis range (Z-values range)
- **G**: Play button and the animation point provide the manual user control over the critical value cut off
- H: Report the appropriate Z-values, density curve height, and cumulative distribution up to Z (i.e., P(X<Z))
- J: There should be a light-colored vertical line at the animation index == Z-value and extending up to the corresponding density height
- K: shaded area represents the integral CDF value .....
- L: Drop-down selected for plotting PDF, CDF or inverse-CDF (quantile) function to plot.

**Starting R Code**: The basic skeleton of one solution (using "**plot\_ly**") is included below. Many solutions are possible and you can start with anything you like, including this initial script.

```
library(magrittr); library(plotly)
select the right user-specified distribution (drop down list)
# Assuming Std Normal N(0,1) going down
# define the range
z < -seq(-4, 4, 0.1)
# points from -4 to 4 in 0.1 steps
# Define the quantile levels for the inverse-CDF (quantile) function)
q<-seq(0.001, 0.999, 0.01)
# probability quantile values from 0.1% to 99.9% in 0.1% steps
# define a DF containing Z, PDF and CDF
dStandardNormal <- data.frame(Z=z, Density=dnorm(z, mean=0, sd=1),
Distribution=pnorm(z, mean=0, sd=1))
# define an index feature
dStandardNormal$ID <- seq.int(nrow(dStandardNormal))</pre>
# Aggregate frames for interactive plot
aggregate by <- function(dataset, feature) {</pre>
feature <- lazyeval::f eval(feature, dataset)</pre>
levels <- plotly:::getLevels(feature)</pre>
aggData <- lapply(seq along(levels), function(x) {</pre>
cbind(dataset[feature %in% levels[seq(1, x)], ], frame = levels[[x]])
})
dplyr::bind rows(aggData)
}
# Apply the aggregate to ID index
dStandardNormal <- dStandardNormal %>% aggregate by (~ID)
# generate the Plot ly object
plotMe <- dStandardNormal \gg plot ly(x = ~Z, y = ~Density, frame = ~frame,
type = 'scatter', mode = 'lines', fill = 'tozeroy', fillcolor="red",
line = list(color = "blue"), text = ~paste("Z: ", Z, "
Density: ", Density, "CDF: ", Distribution), hoverinfo = 'text' ) %>%
layout ( title = "Standard Normal Distribution Distribution",
# Specify the right distribution and its parameters !!!
yaxis = list( title = "N(0,1) Density", range = c(0,0.45),
zeroline = F, tickprefix = "" # density value
),
xaxis = list( title = "Z", range = c(-4,4), zeroline = T, showgrid = T ) ) %>%
animation opts (frame = 100, transition = 1, redraw = FALSE ) %>%
animation slider ( currentvalue = list ( prefix = "Z: " ) )
# display interactive plot
plotMe
```

The optimal solution will include RMD (source) and HTML output (webapp).

# SOCR 2022 MDP Project: Data Science Fundamentals: Spacekime Analytics

SOCR Project Leader: Yueyang Shen, Milen Velev, Ivo Dinov

| Website:          | http://tciu.predictive.space, www.spacekime.org                                  |
|-------------------|--|
| GitHub:           | https://github.com/SOCR/TCIU   |
| Training Modules: | ODE/PDE, Kaluza-Klein Theory (https://en.wikipedia.org/wiki/Kaluza-Klein_theory) |
| GDrive:           | https://drive.google.com/drive/folders/1PMMBR2bzBPubYMpywLkcTkJPyxOKQ4Ag         |

### Description

The SOCR Data Science Fundamentals project will explore new theoretical representation and analytical strategies to understand large and complex data. It will utilize information measures, entropy KL divergence, PDEs, Dirac's bra-ket operators. This fundamentals of data science research project will explore time-complexity and inferential uncertainty in modeling, analysis and interpretation of large, heterogeneous, multi-source, multi-scale, incomplete, incongruent, and longitudinal data.

See The Enigmatic Kime: Time Complexity in Data Science (<u>https://midas.umich.edu/event/midas-seminar-series-presents-ivo-d-dinov-phd-university-of-michigan/</u>) at the University of Michigan Institute for Data Science (MIDAS) Seminar Series, Slidedeck (<u>http://socr.umich.edu/docs/uploads/2018/Dinov\_TCIU\_Kime\_MIDAS\_2018.pdf</u>).

### **Student Skills**

- Physics, math or engineering background
- R programming with RStudio (IDE) experience, and/or Python/Jupyter Notebook, Python

### Project Goals

- Develop, validate and share a new <u>TCIU Python Package</u>, also see <u>TCIU GitHub repo</u>
- Go through the provided materials and references
- Review the current platform (will be provided)
- Perform 3D and 4D Plot\_Ly visualization of complex manifolds, including 5D space-kime and 2D-curved Kime.
- Identify specific case-study and an R&D direction to go deeper into an meaningfully contribute

### Coordinate with team

### Deliverables

- Visualization protocols
- Math proofs of various physics properties in 5D Minkowski spacekime

### **Team Activities**

- Weekly team face-to-face/Zoom meetings
- Code review Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

### Key points

- What is the problem? Use complex-time physics to formulate data science theory & practice
- Why is it important? There is currently no canonical theory for Big Data discovery science
- What is the SOCR Solution? Blend transdisciplinary knowledge to build a new Data Analytic method
- It's real; here it is (in a pilot form) ... demo ... See TCIU Video
- Why should you consider joining this SOCR-MDP Project? High-risk/high-potential yield project.

- Review the websites and listed resources
- TCIU Website: <u>http://tciu.predictive.space/</u>
- TCIU GitHub: <u>https://github.com/SOCR/TCIU/</u>
- Spacekime: <u>www.SpaceKime.org</u>
- The new Data Science/TCIU/Spacekime book (fully accessible) on the <u>publisher</u>, <u>De Gruyter</u>, <u>website</u> (via UMich IP address) and on <u>UM Library</u> (login required) site. More info is available on the <u>TCIU website</u>.

# SOCR 2022 MDP Project: Al for Qualitative-and-Mixed Data

SOCR Project Leaders: Ivo Dinov

 Website:
 ...<TBD>...

 GitHub:
 ...<TBD>...

 Training Modules:
 <see references below>

 GDrive:
 https://drive.google.com/drive/folders/1SGRi8axPaiWIAzYJXwdmSNJ7wXbwt0ay

### Description

Develop new SOCR/DSPA AI/ML tools for representation, modeling, analysis, interrogation, visualization, and interpretation of (1) qualitative data, (2) mixed qualitative-and-quantitative data, and (3) meta-analyses

### **Student Skills**

- EECS, Computer Science (CSE/CS-LSA) and School of Information (SI)
- AI/ML
- UI/UX design, HTML5, JavaScript

### **Project Goals**

- Build and validate AI/ML tools for qualitative data,
- Build and validate AI/ML tools for mixed qualitative-and-quantitative data,
- Build and validate AI/ML tools for meta-analyses

### **Deliverables**

• Stand-alone Rmd notebooks with methods formulation, computational implementation and example utilization using simulated and real (observed) datasets.

### **Team Activities**

- Weekly team meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

# SOCR 2022 MDP Project: Synthetic Image/Volume Generation

SOCR Project Leaders: Yitong, Yongkai, Ivo Dinov GitHub: TBD

 Training Modules:
 see background papers, R-packages, and tutorials listed below

 GDrive:
 https://drive.google.com/drive/folders/1adD9Kp3s\_ZyiIFNCPIA6OmFrfuwA-i0H

### **Description:**

This project aims to build a set of ML/AI tools and complete protocols that provide pragmatic mechanisms to generate realistic 2D images (e.g., faces or animals) or 3D volumes (e.g., brain MRI data). A number of different techniques can be used to generate synthetic data (including images/volumes, e.g., VAE, GAN, diffusion methods, Stochastic corruption, generative, and discriminative models, etc. Start by reading several of the papers listed below and playing with the code provided in them. Then, try some of the 2D, 3D and 4D hyper-volumes (data) referenced below and synthetically generate similar simulated cases.

### **References:**

- Papers
  - <u>https://paperswithcode.com/paper/diffusion-models-beat-gans-on-image-synthesis/review/?hl=32098</u>
  - https://paperswithcode.com/paper/score-matching-model-for-unbounded-data-score-1/review/?hl=38037
  - https://paperswithcode.com/paper/score-matching-model-for-unbounded-data-score-1

### Additional Resources.

- Papers with Code: <a href="https://paperswithcode.com/task/image-generation">https://paperswithcode.com/task/image-generation</a>
- Datasets:
  - Brain Images:
    - Neuroimaging of a large group of healthy individuals from the community (138 subjects), as well as samples of individuals diagnosed with schizophrenia (58), bipolar disorder (49), and ADHD (45): <u>https://openneuro.org/datasets/ds000030/versions/1.0.0</u>. It's easy to get an account and download the data (80GB). All data is de-IDed/anonymized.
  - 2D images:





CIFAR-100



- CelebA
- Fashion-MNIST
- CUB-200-2011
- STL-10

https://www.socr.umich.edu/html/SOCR\_Research.html



- Review these references
  - <u>AI/ML in meta-analyses</u> please see the attached PDF (<u>AI\_MetaAnalyses\_ToolsSummary.pdf</u>), which includes a dozen tools and an extensive article on the topic
  - <u>AI/ML for qualitative data</u>
    - See DSPA <u>Chapter 11 (Association rules)</u> and <u>Chapter 19 (Text Mining/Sentiment analysis)</u>
    - See the examples and tools in this article (AI\_QualitativeData\_SentimentAnalysis.pdf)
    - See the R package RQDA (<u>https://cran.r-project.org/src/contrib/Archive/RQDA/,</u> <u>http://rqda.r-forge.r-project.org/</u>)
  - AI/ML for mixed methods (qualitative-quantitative) data -
    - <u>https://m-clark.github.io/mixed-models-with-R/random\_intercepts.html</u>
    - <u>https://cran.r-project.org/web/packages/boral/index.html</u>
    - <u>http://bbolker.github.io/mixedmodels-misc/glmmFAQ.html</u>
  - Go over <u>DSPA</u> Chapter 22 (Deep Learning): <u>https://www.socr.umich.edu/people/dinov/courses/DSPA\_notes/22\_DeepLearning.html</u>

IDE for development (Eclipse, WebStorm, IntelliJ, Netbeans, ..., RStuio/RMarkdown, Spyder/Py)

# SOCR 2022 MDP Project: VirtualHospital/Synthetic Data

SOCR Project Leaders: Simeone Marino, Johnny Liu, Ronak Shetty, Ivo Dinov

 Website:
 http://gray-rain.com

 GitHub:
 ...

 Training Modules:
 https://www.edx.org/course/health-informatics-data-and-interoperability-stand

 GDrive:
 https://drive.google.com/drive/folders/1hvzoihZMO-6uJIPbObl4rVc\_QDp9i5tD

### Description

This project involves developing a Virtual Hospital (VH) and Synthetic Patient (SP) capability to simulate realistic electronic health records including categorical, discrete, continuous, imaging, text and biomedical specimen data. This relates to the unstructured-<u>DataSifter</u>, synthetic text generation, text-obfuscation, and text-mining/inference.

### Student Skills

- EECS, Stats/Biostats/Math, Computer Science (CSE/CS-LSA) and School of Information (SI)
- UI/UX design, HTML5, JavaScript

### **Project Goals**

- Go through the Training Modules, practice HTML/JS/Angular/Node programming
- Get your GitHub domain going
- Choose 1-2 deliverables, go over current design, start expansion, include unit tests, pilot development
- Coordinate with team

### Deliverables

- HL7/FHIR Interface (see below)
- Learn the current state of the project (Simeone)
- Examine the prior datasets (NHANES, MIMIC, etc.)
- Choose categorical, discrete, continuous, imaging, text and biomedical specimen data and illustrate examples
- Experiment and demo 1,000 VH/SP cases
- Validate VH/SP using machine learning classification, prediction, clustering and analytics modules.

### **Team Activities**

- Weekly team meetings
- Code review (pull/push Github requests)
- Join/present the SOCR All-Hands Calls (twice a month, SOCR Zoom channel)

- Adopt the <u>HL7 XML Format</u>, examine the "<u>Persona</u>" virtual clinical casebooks, and ensure 2 easy pathways (data portals):
- From-Anywhere Into-VirtualHospital, and
- From VirtualHospital Out-to-Anywhere.
- Data Standards and Formats for Information exchange <u>HL7 Format</u> Fast Healthcare Interoperability Resources (<u>FHIR</u>)
- Instrucitons <u>https://fire.ly/2017/10/31/make-your-first-fhir-client-in-r-within-one-hour/</u>
  - R Packages:
    - RonFHIR
    - <u>FHIR/Github</u>, install.packages("remotes"); remotes::install\_github("TPeschel/fhiR")
    - <u>Fhircrackr</u> see the <u>fhircrackr vignettes</u> for many great examples, may need to install the more advanced dev-version: devtools::install\_github("POLAR-fhiR/fhircrackr").

- FHIR is a <u>very useful standard</u> to describe and exchange medical data in an interoperable way.
   FHIR is <u>not useful for statistical analyses of data</u>, since FHIR data is stored in many nested and interlinked resources instead of matrix-like DF structures.
- Use the available public servers, <u>https://hapi.fhir.org/baseR4</u> or <u>http://fhir.hl7.de:8080/baseDstu3</u> as FHIR server endpoint to connect VH clients to.
- Example of executing a <u>FHIR search</u> of the form <u>[base]/[type]?parameter(s)</u>, where <u>[type]</u> refers to the type of resource you are looking for, and <u>[parameter(s)]</u> characterize specific data-search properties:
  - https://hapi.fhir.org/baseR4/Patient?gender=female
- Documentation/Training
  - Read this simple overview or R-based LH7 data XML representation
  - Read this <u>PDF (RonFHIR-Overview-2018-11-15.pdf)</u> .... A must read!
  - Paper (R/<u>RShiny</u>/FHIR): <u>https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5939961/</u>
- Synthetic data generation using Variational Auto-encoders/auto-decoders (AE-AD).
- We can explore data augmentation and synthetic generation variational AE-AD neural networks to generate synthetic cases (using the decoder process, i.e., predict function of the NN model) and explore the similarities between the joint distributions of the original data and the synth-data.
- <u>https://towardsdatascience.com/how-to-generate-new-data-in-machine-learning-with-vae-variational-autoencoder-applied-to-mnist-ca68591acdcf</u>
- <u>https://towardsdatascience.com/auto-encoder-what-is-it-and-what-is-it-used-for-part-1-3e5c6f017726</u>
- https://arxiv.org/ftp/arxiv/papers/1808/1808.06444.pdf
- <u>https://towardsdatascience.com/autoencoders-for-the-compression-of-stock-market-data-28e8c1a2da3e</u>
- SOCR example <u>https://socr.umich.edu/HTML5/ABIDE\_Autoencoder/</u>
- Charlatan R package: <u>https://cran.r-project.org/web/packages/charlatan/index.html</u>
- **synthpop** package provides synthetic –data generation: <u>https://www.r-bloggers.com/generating-</u> synthetic-data-sets-with-synthpop-in-r/ and saving the synthText in diff formats. See this paper.
- stringdist package allows us to compare strings/text: <u>https://cran.r-project.org/web/packages/stringdist</u>
- GAN (generative Adversarial Network) models for synthetic image generation: <u>https://www.r-bloggers.com/conditional-generative-adversarial-network-with-mxnet-r-package/</u>
  - Also see: https://becominghuman.ai/generative-adversarial-networks-for-text-generation-part-1-2b886c8cab10
- See this medical text-GAN mtGAN (Python) paper: <u>https://arxiv.org/pdf/1812.02793.pdf</u>
- Autoencoder approach: <u>https://github.com/stas-semeniuta/textvae</u> and <u>https://www.aclweb.org/anthology/D17-1066.pdf</u>

NEW VH feature – VirtualHospital – BlockchainContracts (Yiwei, Ashish, ...)

- Use the Ether R-package: <u>https://cran.r-project.org/web/packages/ether/index.html</u>
- <u>Goal</u>: VH becomes the mediator between data-governors (owning data) and domain-scientists/dataanalysis (experts). Neither of these 2 stakeholders trust each other (or GR/VH for that matter). However,

using blockchain technology, GR/VH allows the 2 parties to sign contracts using GR VH/DS technology for data-sharing via statistical obfuscation.

- Once a contract is in place and added to the Ether-contract block-chain, the data-governor can use the VH to obfuscate the data (according to the contract agreement terms, initially likely  $\eta = 1$ , for sharing synth data) and share the de-sensitized data with the partner (domain expert).
- Next, the Data-Analyst (expert) can see the data in their GR/VH profile/storage and begin the analytics.
- Finally, the analysts send back reports of their findings to the corresponding data-governors for review
- Based on the results, re-negotiation of contract terms may take place, potentially lowering the obfuscation level ( $0 \le \eta \le 1$ ). This can strengthen partnerships between the stakeholders.

### See these blockchain materials:

https://drive.google.com/drive/folders/1dtQcVKVkJFSKQVPdAboSesQOy2gLBsTg

# SOCR 2022 MDP Project: HTML5/JavaScript

SOCR Project Leaders: Ivo Dinov

Website: various GitHub: https://github.com/SOCR Training Modules: https://github.com/SOCR/socr-tutorials https://drive.google.com/drive/folders/1uZaLGej8NICGfL9NgiFURPvkNIJshQ5F GDrive:

ReadMe File: https://docs.google.com/document/d/1nyJjJqrDq8wRjEjJumScp50q6ns1SL2CoH8p0fAhsFM/

#### I. Kime CircleCloud WebApp index.html

This is a pure JavaScript/HTML5 app that needs some bug-fixing, improvements, and feature enhancements to show the dynamics of the natural attraction-repelling forces.

#### II. SOCR Vase TCIU Model.html

This is an R-source (Rmd) that is knitted/exported as pure HTML5/JS.

#### III. SOCR\_BivariateNormalDistribution Webapp

How to expand the Bivariate-Normal to any Bivariate Distribution? Define the 2 marginal distributions and use this protocol to specify their association (in terms of the correlation) to derive the joint Bivariate distribution PDF:

https://academic.oup.com/mnras/article/406/3/1830/977873

This is an R-source (Rmd) that is knitted/exported as pure HTML5/JS. Below is a starting JavaScript implementing the basic Bivariate Normal Distribution functions.

We want to expand the current app (see it live here and the code (BVN.zip) is here), to include in addition to the MCMC simulation-based (approximate) calculations, which is already implemented, ADD an exact PDF/CDF based calculations. Include a new check-box to allow the user to specify exact vs. approximate calculation (just like we have for with/without using WebGL for the visualization.

```
function normalCDF(X) {
     // using Hastings algorithm with maximal error=10^{-6}
     var T=1/(1+.2316419*Math.abs(X));
     var D=0.3989423*Math.exp(-X*X/2);
     var Prob=D*T*(.3193815+T*(-0.3565638+T*(1.781478+
                T*(-1.821256+T*1.330274))));
     if (X>0) {Prob=1-Prob}
     return Prob
}
function binormalCDF(x,y,R) { // P(X>x,Y>y;R)
   with (Math) {
           var s=(1-normalCDF(x))*(1-normalCDF(y));
           var sqr2pi=sqrt(2*PI);
           var h0=exp(-x*x/2)/sqr2pi;
           var k0=\exp(-y^*y/2)/sqr2pi;
```

```
var h1=-x*h0;
           var k1 = -y * k0;
           var factor=R*R/2;
           s=s+R*h0*k0+factor*h1*k1;
           var n=2;
           while ((n*(1-abs(R))<5)\&\&(n<101)) {
                factor=factor*R/(n+1);
                h2=-x*h1-(n-1)*h0;
                k2=-y*k1-(n-1)*k0;
                s=s+factor*h2*k2;
                h0=h1; k0=k1; h1=h2;
                k1=k2; n=n+1;
           }
           var v=0;
           if (R>.95) {
                v=1-normalCDF(max(h, k))
                s=v+20*(s-v)*(1-R);
           } else if ((R<-.95)&&(h+k<0)) {</pre>
                v=abs(normalCDF(h)-normalCDF(k))
                s=v+20*(s-v)*(1+R);
           }
    }
   return s;
}
function BVN() {
// The following user inputs are necessary
    X,Y, M1, M2 (means), and S1, S2 (sigmal and sigma 2),
     and Rho (correlation)
    Prob="NaN";
    if ((S1<=0)||(S2<=0)){
      alert ("The standard deviations must be positive.")
    } else if ((R<-1)||(R>1)){
      alert("The correlation coefficient must be between -1 and +1.");
    } else {
         h=-(X-M1)/S1;
         k=-(Y-M2)/S2;
         Prob=binormalCDF(h,k,R);
         Prob=Math.round(100000*Prob)/100000;
     }
    return(Prob);
}
```

# SOCR 2022 MDP Project: Java Applets Code → HTML/JS Apps

SOCR Project Leaders: Ivo Dinov

| Website: | various                           |
|----------|-----------------------------------|
| GitHub:  | https://github.com/SOCR           |
|          | https://github.com/SOCR/SOCR-Java |

Training Modules:https://github.com/SOCR/socr-tutorialsGDrive:TBD

### Convert some of the old SOCR Java Applets to modern HTML5/JavaScript apps

Pilot Project: Start with the SOCR Fourier/Wavelet Game applet

SOCR 1D Fourier / Wavelet signal decomposition into magnitudes and phases (Java applet)



Top-panel: original signal (image), white-color curve drawn manually by the user and the reconstructed synthesized (IFT) signal, red-color curve, computed using the user modified magnitudes and phases Bottom-panels: the Fourier analyzed signal (FT) with its magnitudes and phases <u>http://www.socr.ucla.edu/htmls/game/Fourier\_Game.html</u> (Java Applet - Run in Java-enabled IE)

Java app source-code: <a href="https://github.com/SOCR/SOCR-Java/tree/master/src/edu/ucla/stat/SOCR/games">https://github.com/SOCR/SOCR-Java/tree/master/src/edu/ucla/stat/SOCR/games</a>

Additional critical Java Apps to convert to HTML5/JS:

- LLN:
  - http://wiki.stat.ucla.edu/socr/index.php/SOCR\_EduMaterials\_ExperimentsActivities
  - http://socr.ucla.edu/htmls/exp/LLN\_Simple\_Experiment.html
- CLT
  - <u>http://wiki.stat.ucla.edu/socr/index.php/SOCR\_EduMaterials\_ExperimentsActivities</u>
  - <u>http://socr.ucla.edu/htmls/exp/Sampling\_Distribution\_CLT\_Experiment.html</u>
- $\pi$  and *e* stochastic estimation experiments:
  - http://socr.ucla.edu/htmls/exp/Uniform\_E-Estimate\_Experiment.html
- Polynomial model fitting:
  - <u>http://socr.ucla.edu/Applets.dir/SOCRCurveFitter.html</u>

## SOCR 2022 MDP Project: <u>Topological Data Analysis (TDA),</u> Persistent Homology, and Betti Numbers for Point-cloud Data

SOCR Project Leaders: Ribhu (Antareeksh Deb), Ivo DinovWebsite:TBDGitHub:TBD

 Training Modules:
 see background papers, R-packages, and tutorials listed below

 GDrive:
 https://drive.google.com/drive/folders/1QZLjH\_xN\_SSsrjF2jsc9mDQBUCdDKPIY

### Description:

This project aims to expand the SOCR lab data analytical capabilities using advanced topological representations. Data are viewed as high-dimensional point clouds. These are interpreted as samples through a high-dimensional manifold which will be modeled using simplicial complexes. Using the simplicial decomposition, we can compute persistent homology, Betti numbers, and derive other topological metrics that can be used for classification, regression and clustering.

### **References**:

- Papers
  - <u>https://www.ams.org/journals/notices/201905/rnoti-p686.pdf</u>
  - https://arxiv.org/pdf/1705.02037.pdf
  - https://arxiv.org/pdf/1812.02987.pdf
- Tutorials
  - <u>See DSPA Chapter 5</u>:
    - http://www.socr.umich.edu/people/dinov/courses/DSPA\_notes/05\_DimensionalityReduction.html
  - http://www.stat.cmu.edu/~jisuk/files/20180613\_SoCG\_Jisu\_KIM\_TDA\_slide.pdf
- R-Code
  - https://cran.r-project.org/web/packages/TDA/TDA.pdf
  - https://cran.r-project.org/web/packages/TDAstats/TDAstats.pdf
- Test data: https://umich.instructure.com/courses/38100/files/folder/Case\_Studies

### Topological, Fiber-Bundles, and Differential-Geometric approaches to Data Science

- Flag Manifolds, see <u>nested Flag vector spaces and Canonical Correlation Analysis (CCA) (DS applications</u>)
- Grassmann Manifolds, see; Foundations of Grassmann manifold
- <u>Comparing datasets using flag-manifolds</u>.

### Algebraic Data Analysis (ADA)

- Similarly to TDA, try to design a new Dataset → Algebraic Group, Ring or Field mapping that transforms each dataset into a mathematical object (e.g., a Lie group like SO(n,F) with a corresponding Lie algebra so(n,F) …
- See:
  - https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=9259191
  - http://www.cs.technion.ac.il/~ron/PAPERS/Conference/RosBroBroKim3DOR2012.pdf
  - https://arxiv.org/pdf/1912.00396.pdf

# SOCR 2022 MDP Project: Novel AI/ML Techniques

SOCR Project Leaders: Ivo Dinov Website: TBD

GitHub: TBD

 Training Modules:
 see background papers, R-packages, and tutorials listed below

 GDrive:
 https://drive.google.com/drive/folders/1CSg6xxcwYsjRSV1AjgFNjm1dfTQ5jxeS

### **Description:**

This project aims to build new <u>reinforcement-learning</u>-based and <u>deep neural network</u>-based techniques for representation of complex high-dimensional data and for clustering, classification, interpolation, extrapolation, forecasting, and prediction.

### References:

- Papers
  - See Project partition:
    - https://drive.google.com/drive/folders/1CSg6xxcwYsjRSV1AjgFNjm1dfTQ5jxeS
  - See this talk: <u>https://www.youtube.com/watch?v=8QLDXITiRII</u>
  - https://h2o.gitbooks.io/h2o-tutorials/content/tutorials/ensembles-stacking/ (DNN + SuperLearner)
- Tutorials
  - <u>See DSPA Appendix: https://dspa.predictive.space</u>
- R-Code
- o ...
- Test data: <u>https://umich.instructure.com/courses/38100/files/folder/Case\_Studies</u>
- Review these concepts:
  - <u>https://en.wikipedia.org/wiki/Backpropagation</u>
  - <u>https://en.wikipedia.org/wiki/Invariant\_subspace</u>
  - Learn to compress and compress to learn
  - DSPA: dspa.predictive.space

**Demos**: Here are some specific examples we can consider implementing using QL in the Appendix:

- Puzzle: <u>https://en.wikipedia.org/wiki/15\_puzzle</u>
- <u>Non-linear function optimization</u> like this DSPA example: https://socr.umich.edu/people/dinov/courses/DSPA\_notes/21\_FunctionOptimization.html#93\_Convexity
- <u>Games</u>:
  - 2048: <u>https://towardsdatascience.com/a-puzzle-for-ai-eb7a3cb8e599</u> and <u>https://github.com/voice32/2048\_RL</u>
  - Snake: https://towardsdatascience.com/how-to-teach-an-ai-to-play-games-deep-reinforcementlearning-28f9b920440a
  - Atari: <u>https://becominghuman.ai/lets-build-an-atari-ai-part-1-dqn-df57e8ff3b26</u> and datasets (<u>https://paperswithcode.com/task/atari-games</u>)
- Some harder problems are included here: <u>https://analyticsindiamag.com/reinforcement-learning-top-state-of-the-art-games-alphago/</u>

# SOCR 2022 MDP Project: Synthetic Image/Volume Generation

SOCR Project Leaders: Ivo Dinov GitHub: TBD

**Training Modules**: see background papers, R-packages, and tutorials listed below **GDrive**: https://drive.google.com/drive/folders/1adD9Kp3s\_ZyiIFNCPIA6OmFrfuwA-i0H

### **Description:**

This project aims to build a set of ML/AI tools and complete protocols that provide pragmatic mechanisms to generate realistic 2D images (e.g., faces or animals) or 3D volumes (e.g., brain MRI data). A number of different techniques can be used to generate synthetic data (including images/volumes, e.g., VAE, GAN, diffusion methods, Stochastic corruption, generative, and discriminative models, etc. Start by reading several of the papers listed below and playing with the code provided in them. Then, try some of the 2D, 3D and 4D hyper-volumes (data) referenced below and synthetically generate similar simulated cases.

### **References:**

- Papers
  - <u>https://paperswithcode.com/paper/diffusion-models-beat-gans-on-image-synthesis/review/?hl=32098</u>
  - https://paperswithcode.com/paper/score-matching-model-for-unbounded-data-score-1/review/?hl=38037
  - https://paperswithcode.com/paper/score-matching-model-for-unbounded-data-score-1

### Additional Resources.

- ٠
  - See <u>DSPA Chapter 22 (Deep Learning)</u>
- Papers with Code: <a href="https://paperswithcode.com/task/image-generation">https://paperswithcode.com/task/image-generation</a>
- Datasets:
  - Brain Images:
    - Neuroimaging of a large group of healthy individuals from the community (138 subjects), as well as samples of individuals diagnosed with schizophrenia (58), bipolar disorder (49), and ADHD (45): <u>https://openneuro.org/datasets/ds000030/versions/1.0.0</u>. It's easy to get an account and download the data (80GB). All data is de-IDed/anonymized.
  - 2D images:
    - CIFAR-10
    - ImageNet
    - MNIST
    - CIFAR-100
    - Cityscapes
    - CelebA
    - Eachian MN
    - Fashion-MNIST
    - CUB-200-2011
    - STL-10
    - Oxford 102 Flower
    - <u>100\*100 grid/table of 2D PNG brain (MRI) images</u>