Hands-on Training - Data Science Challenge

Simon McIntosh ITER Organization

Nathan Cummings Culham Center for Fusion Energy

Thursday 22nd May 2025

Joint ICTP-IAEA Fusion Energy School Trieste, Italy





This presentation focuses on applied Data Science for Fusion.

The elephant in the room.



Applied Data Science in Fusion.





Challenge Notebooks.











The elephant in the room is over-fitting.



Hands-on Training - Data Science Challenge Simon McIntosh and Nathan Cummings 22/05/2025



"All models are wrong, but some are useful". George Box Non-linear parametrizations can fit arbitrarily complexity structures.



Hands-on Training - Data Science Challenge Simon McIntosh and Nathan Cummings 22/05/2025



"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk".



Enrico Fermi



John von Neumann

Hands-on Training - Data Science Challenge Simon McIntosh and Nathan Cummings 22/05/2025



Simple basis functions can generate complex shapes.



Scientific Data is a valuable product of expensive experiments. If not properly archived the value of this data depreciates rapidly.



Capital cost ~0.5 billion 2014 US dollars



Operating cost ~200,000 euros per day



Code Data is read more often than it is written. Code Data should always be written in a way that promotes readability.

Data Science Challenges of the ITER International School 2024



MAST Plasma Current

Infer plasma current produced by CCFE's Mega Ampere Spherical Tokamak from discrete magnetic diagnostic data. https://github.com/Simon-McIntosh/data-sciencechallenges





MAST Plasma Volume

Infer the volume of plasmas produced by the CCFE's Mega Ampere Spherical Tokamak using frames from a wide-angle visible spectrum camera.



MAST Plasma Equilibrium

Infer two-dimensional poloidal flux maps produced by the EFIT++ equilibrium reconstruction code from a diverse set of diagnostic measurements.

Challenge #1 MAST Plasma Current



Challenge #2 MAST Plasma Volume



Challenge #3 MAST Plasma Equilibrium

Data variables:	
center_column	<pre>(center_column_channel, time) float64 16kB</pre>
coil_currents	<pre>(coil_currents_channel, time) float64 19kB</pre>
coil_voltages	<pre>(coil_voltages_channel, time) float64 13kB</pre>
flux_loops	(flux_loops_channel, time) float64 19kB
outer_discrete	(outer_discrete_channel, time) float64 26kB .
saddle_coils	<pre>(saddle_coils_channel, time) float64 13kB</pre>
dalpha_mid_plane_center	(time) float64 3kB
dalpha_mid_plane_wide	(time) float64 3kB
dalpha_tangential	(time) float64 3kB
hcam_l	(hcam_l_channel, time) float64 58kB
hcam_u	(hcam_u_channel, time) float64 58kB
ne	(time, major_radius) float64 210kB
ne_core	(time) float64 3kB
pe	(time, major_radius) float64 210kB
te	(time, major_radius) float64 210kB
te_core	(time) float64 3kB
<pre>shot_index</pre>	(time) float64 3kB
magnetic_flux	(time, z, major_radius) float64 14MB
tcam	(tcam_channel, time) float64 58kB



In summary this talk has warned you of the dangers of overfitting and has given you the opportunity to learn more via the challenges.

Remember the elephant.

Data Science Challenge facilitated by FAIR data and open-source tools.

Image: state of the state of the

Practical application is considered the best way to learn.





"Tell me and I forget, teach me and I remember, involve me and I learn."

Hands-on Training - Data Science Challenge Simon McIntosh and Nathan Cummings 22/05/2025

