

Queen Mary **University of London**



DISCnet Placement at APSensing UK

-DAS Prognostic Health Monitoring with Machine Learning-



Arran Freegard Supervisors: U. Blumenschein, A. Belyaev & S. Moretti





APSensing UK

- HP (1939) and Agilent Technologies (1999) heritage
- German based optical sensing technologies distributor (for over 25 years):
 - Distributed Temperature Sensors (DTS)
 - Distributed Acoustic Sensors (DAS)
 - Distributed Thermal Gradient Sensors (DTGS)
- Subsea and land cables with DAS and DTS power cable monitoring systems
- UK based Data Science Office in Basingstoke, focusing on Alarm systems using Machine Learning techniques



Project I: Telemetry Dashboards with ML predictions

- Analyse time series data produced by DAS systems to monitor hardware health
- Parameters recorded and sent to a database periodically:
 - Disk and internal Temperatures
 - Fan Speeds
 - Disk Usages
 - Status/Error codes
 - Uptime/Reboot cycles
- Build a predictive model to determine current operating health status and future failure

					Inte	erna	alT	en	np	1							diskl	Jsag	ge1,2,	3	con	stant			•
/ar) SI	d p	ISIN	I G							(dis	kΤ	em /	p1	2,3	3 F	- Fan1	.2	DA	S Syste	ems wit	n Telem	etry	
	Status			Serial		Internal (degC	lemp C)	Dis	k Tem	p (deg	C)	Total	Disk Usa	age (TB)	Fan (R	Speed PM)		Machine Lea	arning Verdict		Reboot	Days	Total days		
DAS	IU P	J	Date / Time	Number	Uptime	IU	PU	nvme	sda	sdb	sdc	sda	sdb	sdc	Fan1	Fan2	Linear Regressions	<mark>SVM</mark>	IsolationForest	Autoencoder	Cycles	Booted	booted	Version	Update Date
0	1 LIV	E	16/12/19 23:30	DE5000PP02	1w 14h 36m	52	48	48	44	45	44	918.7	918.5	53.2	5900	5900	Normal	Normal	Normal	Normal	579	167	231	1.1.4.3643	09/12/2019
1	1 Star	dby	16/12/19 23:30	DE52000218	4d 13h 10m	50	51	45	44	44	44	8.6	8.6	1.7	6200	6190	Normal	Normal	Normal	Normal	25	4	74	1.1.4.3644	09/12/2019
0	1	E	16/12/19 23:30	DE52000235	3w 3d 14h 48m	45	52	52	46	46	46	10.7	10.7	0.4	6800	6800	Normal	ANOMALY	Normal	Normal	27	24	30	1.1.4.3526	03/12/2019
0	1 LIV	E	16/12/19 23:30	DE52000250	6w 5d 9h 48m	33	44	40	39	39	38	8.6	8.6	1.9	5100	5200	Normal	Normal	Normal	Normal	19	46	117	1.1.4.3213	12/11/2019
0	1 LIV	E	16/12/19 23:30	DE52000LP2	1w <mark>4d 1</mark> 6h 10m	37	39	46	44	44	42	398.6	398.5	54.7	5800	5900	Normal	Normal	Normal	Normal	58	11	149	1.1.3.3094	19/11/2019
1	0	-NOT	16/12/19 23:30	DE5200DV01	7w 4d 12h 34m	0	43	35	36	36	36	6.1	6.1	0.1	4900	4900	Normal	ANOMALY	Normal	Normal	13	200	200	1.1.4.3423	25/11/2019
	ō	-N(15/12/19 23:36	DE52000LP1	2d 9h 12m	ERROR	42	40	36	37	37	277.5	277.5	277.5	5200	5100	FanRMSE: +/- 1759.33	Normal	ANOMALY	ANOMALY	6	2	188	1.1.3.3094	11/11/2019
	1 0	¢	04/12/19 11:33	DE52000219	34m	0	28	24	26	31	28	65.4	58.1	61.9	3800	3600	nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY	75	0	32	1.1.0	21/06/2019
1	1 0	¢	28/11/19 11:30	DE52000LP3	6d 17h 42m	46	41	40	37	38	37	0.0	0.0	0.0	4100	3900	Normal	Normal	ANOMALY	Normal	64	12	177	1.1.4.3423	27/11/2019
	1 0	¢	24/11/19 23:30	DE52000253	1w 3d 6h 35m	37	32	34	28	28	29	0.0	0.0	0.0	3900	3800	Normal	Normal	ANOMALY	ANOMALY	28	10	10	1.1.3.3094	21/11/2019
1	1 0	<	15/09/19 23:31	DE5200PP04	6d 12h 35m	42	40	44	42	42	40	175.9	175.9	9.7	5100	5000	Normal	Normal	ANOMALY	Normal	104	6	82	1.1.1.2584	30/08/2019
1	0 ive: -IU	NOT r	19/08/19 23:30	DE52000PP1	5d 9h 23m	0	45	40	41	41	40	8.0	8.0	1.9	5400	5500	FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal	18	5	26	1.1.1.2245	24/07/2019
1	1	-NO	23/07/19 23:30	DE52000201	5w 5d 6h 41m	ERROR	46	40	41	42	41	8.0	8.0	1.9	5300	5400	Normal	Normal	ANOMALY	Normal	16	39	105	1.1.0	11/06/2019
	1	-NO	03/07/19 23:31	DE52000229	12h 31m	ERROR	45	39	41	42	41	0.7	0.7	0.1	4000	4100	Normal	Normal	ANOMALY	Normal	29	0	0	1.1.0	04/06/2019
•		ready	29/05/19 12:50	DE5xxxxxxx	4w 5d 1h 59m	0	38	28	34	35	34	1.1	1.1	0,1	4400	4400	Normal	Normal	ANOMALY	Normal	0	75	87	1.0.10	0.0

Telemetry Dashboard #1: Variables

Machine learning prediction's based on different models

Telemetry Dashboard #1: Manual Thresholds



Red text: Has an error

IU not present, but files left over

Telemetry Dashboard #1: Machine Learning Columns

ML Process:

- Trim data to remove anomalies for training:
 - Certain TBD12 codes
 - No command timeouts
 - Systems that have been running for more than a day
 - Temperatures less than 60
 - Fan speeds less than 7000
 - Removed data 12 hours before shut downs
- Trains (and validate) models on this healthy data set
- diskUsage1, diskUsage3, diskTemp1, Fan1, internalTemp1, TBD2 (PUTemp), uptime, BootTimer, Live (based on VIF analysis)
- NOT trained per ID (this gave bad results, not enough data)
- Append new data and train, if not anomalous
- 80-90% validation
- Train on "Healthy" data 12 hours before to test new incoming data

Machine Learning verdict											
Linear Regressions	SVM	IsolationForest	Autoencode								
Normal	Normal	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	ANOMALY	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	ANOMALY	Normal	Normal								
FanRMSE: +/- 1759.33	Normal	ANOMALY	ANOMALY								
nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	ANOMALY								
Normal	Normal	ANOMALY	Normal								
FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
FanRMSE: +/- 1260793.473	Normal	ANOMALY	Normal								

Data Preprocessing

Scalers

- Re-scale data (eg: normalized, min-max scaling)
- Some models train better on different methods of scaled data, giving better results

Principal Component Analysis (PCA)

- Reduce dimensionality of training set by analysing relations between variables
- Determining which variables skew the data the most
- Can make ML learning more efficient and run faster
- Use Linear regressions to find Variance Inflation Factor (VIF) and remove components of VIF>5

diskUsage1			0.02	0.70	0.0		0.57	0.54	0.52	0.49	0.53		-0.13	-0.55	0.51	0.55	0.55	0.55		
diskUsage2 -							0.57	0.54	0.52	0.49	0.53		-0.13		0.51	0.99				
diskUsage3 -							0.45	0.41	0.42	0.42	0.42	0.54	-0.087		0.41	0.8				- 0.
diskTemp1 -				1							0.9		-0.055		0.58					
diskTemp2 -													-0.064							
diskTemp3 -													-0.049		0.57					- 0.
Fanl -	0.57	0.57	0.45										-0.016		0.58					
Fan2 -	0.54	0.54	0.41										-0.031		0.58	0.57	0.57	0.57		
internalTemp1 -	0.52	0.52	0.42										0.044		0.49	0.54	0.54	0.54		
internalTemp2 -	0.49	0.49	0.42										0.013	-0.47	0.47	0.51	0.51	0.51		- 0.
TBD2 -	0.53	0.53	0.42										0.051		0.5	0.56	0.56	0.56		
TBD5 -			0.54										-0.018		0.44	0.76				
TBD12 -	-0.13	-0.13	-0.087	-0.055	-0.064	-0.049	-0.016	-0.031	0.044	0.013	0.051	-0.018	1	0.17	-0.17	-0.15	-0.15	-0.15		(
Status -										-0.47			0.17			-1				
uptime -	0.51	0.51	0.41	0.58		0.57	0.58	0.58	0.49	0.47	0.5	0.44	-0.17			0.55	0.55	0.55		
Cycle_count1 -								0.57	0.54	0.51	0.56		-0.15		0.55	1				(
Cycle_count2 -								0.57	0.54	0.51	0.56		-0.15		0.55	1				
Cycle_count3 -	diskUsage1	diskUsage2	diskUsage3 –	diskTemp1 -0	diskTemp2 -6	diskTemp3 - 0	Fanl o	0.57 Eau5	internalTemp1 60	1400 InternalTemp2	0.56	0.76 1802	15.0.	Status	0.55 nhtime	Cycle_count1 -	Cycle_count2 -	Cycle_count3 -		

IU Temp TBD2	- 0.53	0.53	0.42	0.9	0.89	0.9	0.94	0.93	0.95	0.91	1	0.83	0.051	-0.53	0.5	0.56	0.56	0.56
PU Temp TBD5	- 0.75	0.75	0.54	0.92	0.91	0.91	0.79	0.78	0.83	0.78	0.83	1	-0.018	-0.72	0.44	0.76	0.76	0.76
IU Status TBD12	0.13	-0.13	-0.087	-0.055	-0.064	-0.049	-0.016	-0.031	0.044	0.013	0.051	-0.018	1	0.17	-0.17	-0.15	-0.15	-0.15
Status	0.99	-0.99	-0.8	-0.76	-0.79	-0.75	-0.58	-0.55	-0.51	-0.47	-0.53	-0.72	0.17	1	-0.57	-1	-1	-1
uptime	- 0.51	0.51	0.41	0.58	0.59	0.57	0.58	0.58	0.49	0.47	0.5	0.44	-0.17	-0.57	1	0.55	0.55	0.55
Cycle_count1	0.99	0.99	0.8	0.78	0.81	0.78	0.6	0.57	0.54	0.51	0.56	0.76	-0.15	4	0.55	1	1	1
Cycle_count2	- 0.99	0.99	0.8	0.78	0.81	0.78	0.6	0.57	0.54	0.51	0.56	0.76	-0.15	4	0.55	1	1	1
Cycle_count3	diskUsage1 - 5	diskUsage2	diskUsage3 –	diskTemp1 -	diskTemp2 -0	diskTemp3	Fanl J	Fan2 0	nternalTemp1 6	nternalTemp2 4	1BD2 0	0.76 EDD2	BD12 Ó	Status	0.55 nbtime	Cycle_count1 . =	Cycle_count2 .=	Cycle_count3

Data Preprocessing: Linear Regression Analysis



Data Preprocessing: Linear Regression Analysis



Data Preprocessing: Linear Regression Analysis



Telemetry Dashboard #1: Machine Learning Columns

	Linear Regressions	SVM	IsolationForest	Autoencoder
	Normal	Normal	Normal	Normal
	Normal	Normal	Normal	Normal
Linear Regressions:	Normal	ANOMALY	Normal	Normal
 Skiearn Linear Regression model No scaling or PCA 	Normal	Normal	Normal	Normal
Checks predictions for	Normal	Normal	Normal	Normal
InternalTemp1	Normal	ANOMALY	Normal	Normal
• Fan1	FanRMSE: +/- 1759.33	Normal	ANOMALY	ANOMALY
PU Temp	nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY
	Normal	Normal	ANOMALY	Normal
	Normal	Normal	ANOMALY	ANOMALY
 Flags if RMSE > 10% of actual value, 	Normal	Normal	ANOMALY	Normal
per three variables	FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal
	Normal	Normal	ANOMALY	Normal
	Normal	Normal	ANOMALY	Normal
	Normal	Normal	ANOMALY	Normal
	FanRMSE: +/- 1260793.473	Normal	ANOMALY	Normal

Support Vector Machine (SVM)

- SKLearn ML model that takes labeled training data
- Outputs an "optimal hyperplane", categorizes new examples

1 Class SVM

- Unsupervised outlier detection
- Estimates the support of a high-dimensional distribution
- Trains on "Healthy" class of data
- Tests new input data to determine if it fits within criteria of trained data
- Train on this new data if it is determined to be healthy, updates model

Telemetry Dashboard #1: Machine Learning Columns

Machine Learning Verdict									
Linear Regressions	SVM	IsolationForest	Autoencoder						
Normal	Normal	Normal	Normal						
Normal	Normal	Normal	Normal						
Normal	ANOMALY	Normal	Normal						
Normal	Normal	Normal	Normal						
Normal	Normal	Normal	Normal						
Normal	ANOMALY	Normal	Normal						
FanRMSE: +/- 1759.33	Normal	ANOMALY	ANOMALY						
nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY						
Normal	Normal	ANOMALY	Normal						
Normal	Normal	ANOMALY	ANOMALY						
Normal	Normal	ANOMALY	Normal						
FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal						
Normal	Normal	ANOMALY	Normal						
Normal	Normal	ANOMALY	Normal						
Normal	Normal	ANOMALY	Normal						
FanRMSE: +/- 1260793.473	Normal	ANOMALY	Normal						

1 Class SVM:

- min_max_scaler
- PCA for 2 components
- Nu parameter: 0.1
- gamma='auto'
- Anomaly if SVM returns -1

Isolation Forest

- SKLearn ML anomaly detection model
- Isolates observations by randomly selecting a feature
 - then randomly selecting a split value between max and min values of that feature
- "Recursive partitioning" represented by a tree structure:



- > Number of splittings required to isolate a sample == **path length** from the root node to the terminating node
- > **Path length**, averaged over a forest of random trees == measure of normality and decision function
- Random partitioning produces noticeably shorter paths for anomalies
 - > If forest of random trees all produce shorter path lengths for samples, are likely to be anomalies

Telemetry Dashboard #1: Machine Learning Columns

Linear Regressions	SVM	IsolationForest	Autoencoder								
Normal	Normal	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	ANOMALY	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	Normal	Normal	Normal								
Normal	ANOMALY	Normal	Normal								
FanRMSE: +/- 1759.33	Normal	ANOMALY	ANOMALY								
nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	ANOMALY								
Normal	Normal	ANOMALY	Normal								
FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
Normal	Normal	ANOMALY	Normal								
FanRMSE: +/- 1260793.473	Normal	ANOMALY	Normal								

Machine Learning Verdict

IsolationForest:

- min_max_scalerPCA for 2 components

Anomaly if returns -1

Autoencoders

- Artificial neural network
- Unsupervised ML model that learns from input data
- Given unlabeled training examples
- Applies backpropagation, sets target values to be equal to the inputs
- Gain insight on structure of the data by applying constraints to the network: limit number of hidden layers



Autoencoders

- For small number of hidden units, network forced to learn a "compressed" representation of the input, must try to "'reconstruct" the input
- For random input, compression task would be very difficult.
- But if there is structure in data, eg: correlated input features, then this algorithm will be able to discover some of those correlations.
- Simple autoencoder learns a low-dimensional representation similar to PCAs.
- "Overkill" for simple data, can lead to over-fitting



Telemetry Dashboard #1: Machine Learning Columns

Machine Learning Verdict

Linear Regressions	SVM	IsolationForest	Autoencoder
Normal	Normal	Normal	Normal
Normal	Normal	Normal	Normal
Normal	ANOMALY	Normal	Normal
Normal	Normal	Normal	Normal
Normal	Normal	Normal	Normal
Normal	ANOMALY	Normal	Normal
FanRMSE: +/- 1759.33	Normal	ANOMALY	ANGMALY
nvme Temp RMSE: +/- 14755.6	Normal	ANOMALY	ANOMALY
Normal	Normal	ANOMALY	Normal
Normal	Normal	ANOMALY	ANOMALY
Normal	Normal	ANOMALY	Normal
FanRMSE: +/- 3552.727	Normal	ANOMALY	Normal
Normal	Normal	ANOMALY	Normal
Normal	Normal	ANOMALY	Normal
Normal	Normal	ANOMALY	Normal
FanRMSE: +/- 1260793.473	Normal	ANOMALY	Normal

Autoencoder:

- min_max_scaler
- No PCA
- 9 initial layers, 1 per feature

min_max_scaler.fit(train_X)
Apply transform to both the training set and the test set.

train_X = min_max_scaler.transform(train_X)
test_X = min_max_scaler.transform(test_X)

No of Neurons in each Layer [9,6,3,2,3,6,9] input_dim = train_X.shape[1] encoding_dim = 6 input_layer = Input(shape=(input_dim,)) encoder = Dense(encoding_dim, activation="tanh", activity_regularizer=regularizers.ll(10e-5))(input_layer) encoder = Dense(int(encoding_dim / 2), activation="tanh")(encoder) encoder = Dense(int(encoding_dim / 2), activation="tanh")(encoder) decoder = Dense(int(encoding_dim / 2), activation="tanh")(encoder) decoder = Dense(int(encoding_dim / 2), activation='tanh')(decoder) decoder = Dense(int(encoding_dim, activation='tanh')(decoder) decoder = Dense(int(encoding_dim, activation='tanh')(decoder) autoencoder = Model(inputs=input_layer, outputs=decoder) autoencoder.compile(optimizer='adam', loss='mse') history = autoencoder.fit(train_X, train_X,epochs=nb_epoch,batch_size=batch_size,shuffle=True,validation_split=0.1,verbose=0) df_results['AAMSE']=np.mean(np.power(test_X - autoencoder.predict(test_X), 2), axis=1)

Anomaly if MSE>0.1



Telemetry Dashboard #2:





Project II: Machine Learning Library Upload and Dashboard

- DAS systems alarm based on event types: Trains, Cars, Cable thefts
- Different ML libraries built to classify events
- Visualizations of events and model performance from updating database
 - Automated and uploaded to dashboard

Machine Learning Library: Upload Data Script

- 'Upload_Data.py', gives hdf5 file location
- Takes features from Feature_Calculation file, feature_sizes
- Class from Folder names
- Sitename based on "Metadata/Country"
- Generates hash id, based on file name and row number, allows for overwriting on second upload
- Row is just index
- If file unspecified, will run on all ML library files



Machine Learning Library: Dashboard



- 'plots.py' running on Venus, pulls Elasticsearch data makes plots for Dashboard, curl plots to Pluto
- 'script.py' running on **Pluto** (/usr/lib/cgi-bin/ml_dashboard/, ith plots read from /var/www/html/ml_plots simlink to home/ml_plots)







Average RMSE per Class





Average PeaktoPeak per Class







Average FFT



pit.step(x,womp2p, label+classtype, where="mid",color="c"+str((chane.index(classtype)



- Train

- Ship

- Car

Ship Background

Digging Mechanical

Prysmian_Digging

- Track Distance

---- Prysmian_Train

---- CableTheft

Orsted

Drilling

Dropping

---- RTE_Version7_Cars

Digging Manual

175

Aggriculture Machine

---- Walking

--- Boring



Summary and Outcomes

- Explored real implementations of ML models for data science, in a workplace scenario:
 - Successfully implemented ML model to **detect anomalies** in incoming data (see plot)
 - Created **updating databases** for new **ML statistics** and **visualizations**



Summary and Outcomes

Improvements of Data Science Skills:

- Data exploration, wrangling and analysis, understanding and filtering valuable information
- Appropriate usage of different **ML models** (Linear regression, 1 class svm, autoencoders)
- Various **programming** languages and modules

Improvements of Soft Skills:

- Project time management, task priority assessing and keeping to deadlines
- Liaising/teamwork with colleagues on project details/plans
- Taking **leadership** on a project and reporting findings
- Tackling issues relating to project difficulties/compromises (small amounts of data, limits to software packages, etc.)
- Group meetings and scheduling
- Creative thinking and approaches to different data types
- **Presenting** results and project hand-off details/instructions