# SageMaker ML Algorithms

Wednesday, 27 July 2022    14:33

| | Algorithm | Input | Usage | Hyperparameter | Instance Types |
|---|---|---|---|---|---|
| Linear Learner | - **Linear regression - numeric quantity**<br>- **Linear Threshold Fun: binary and multi-class categorical classification**<br>- **Linear relationship**: x,y plot of feature and target = staright line diagonaly<br>- Property price - number of rooms | - **RecordIO-wrapped protobuf (Float32)**<br>- CSV (first column as label),<br>- **Reader**: File or Pipe (performant as it stream from S3) mode | - **Pre-process:** Normalized and shuffled data<br>- **Training**: stochastic gradient<br>- **Optimization Algo**: SGD, Adam,etc<br>- **Tune**: L1 & L2 regularization | - Balance_multiclass_weigths<br>- Learning_rate<br>- min_batch_size<br>- Regularization (L1)<br>- Weight Decay (L2) | - **Single or multi-machine CPU**<br>- **Single GPU** |
| XGBoost | - **eXtreme Gradient boosting:** Boosted group of **decision trees**<br>- **Corrects Decision trees errors from previous model to new model**<br>- **Classification**<br>- **Regression trees** | - CSV or **libsvm**<br>- AWS: recordIO-protobuf & Parquet | - **Model**: Serialize/Deserialize models with Pickle<br>- **Framework in SM :** Sagemaker.xgboost<br>- **Built-in algorithm** | - **Prevent overfitting**: Subsample, Eta<br>- **Step size shrink**: Eta<br>- **Minimum loss reduction to create partition**: Gamma<br>- **L1** : Alpha<br>- **L2**: Lambda<br>- **Optimize on AUC (area under curve), error, rmse**: Eval_metric<br>- **Adjust positive and negative weights**: scale_pos_weight<br>- **Max depth of tree**: max_depth | - **Multi-instance CPU**<br>- **Memory-bound**: M5<br>- **XGBoost 1.2 Single GPU**: P3, tree_method_hyperparameter=gpu_hist **(Quick and cost effective)** |
| Seq2Seq (DL) | - **RNN's and CNN's** with attention<br>- **Machine translation**<br>- **Text summarization**<br>- **Speech to Text** | - **tokenized text files**<br>- **Convert to recordIO-protobuf (Integer** tokens)<br>- training data, validation data<br>- **vocabulary files** | - **Pre-trained models** in SM<br>- Can take days to train new<br>- **Public training dataset for translation** | - Batch_size<br>- **Optimizer_type** = adam, sgd<br>- Learning_rate<br>- Num_layers_encoder<br>- Num_layers_decoder<br>- Optimize on:<br>  • **Accuracy**<br>  • **BLUE & Preplexity:** Machine translations | - **Only GPU** (P3) as DL Algo<br>- **Single machine with multi-GPU** |
| DeepAR (DL) | - **RNN's**<br>- **Forecast one dimensional time-series data**<br>- Train same model over related time series<br>- **Frequencies and Seasonality** | - **JSON lines** in file<br>- **Gzip** or **Parquet**<br>- **Record MUST**: **Start timestamp, target values**<br>- **Record CAN**: **Dynamic**_features, **Categorical** feat | - **entire time series for Training, test and inference**<br>- remove last timepoints from training<br>- **train on many time series**<br>- **Prediction length < 400** | - **Context_length**: Number of time points before making prediction<br>- Epochs<br>- Mini_batch_size<br>- Learning_rate<br>- Num_cells | - **Single or multi CPU**<br>- **Single or multi GPU**<br>- Start with CPU<br>- **Inference**: **CPU**<br>- **Tuning**: **Larger** instance |
| BlazingText (DL) | - **Supervised Text classification**: Label sentences, NOT DOCs<br>- **Word2vector (word embedding)**: Machine **translation, sentiment analysis** | - **Text classification**: One sentence per line, first **word "_label_ followed by label"**<br>- **Word2vec**: text file with one sentence per line | - **Word2Vec input modes**: Cbow, Skkip-gram, Batch-skip gram (computation over many CPU nodes) | - **Word2vec:**<br>  - Mode<br>  - Learning_rate, window_size, vector_dim, negative_samples<br>- **Text Classification:**<br>  - Epochs, Learning_rate, Word_grams, Vector_dim | - **Cbow & skinpgram:** Single CPU or GPU (P3)<br>- **Batch_skipgram:** Single or multiple CPU<br>- **Text classification:** CPU if training data < 2 GB otheriwse GPU (P2, P3) |
| Object2Vec (DL) | - **Object embedding, CNN**<br>- Compute nearest neighbours of objects<br>- **Recommender (similar items)**<br>- Can be used for entire Docs instead of just objects<br>- **Ex: sentence Embedding** | - **Tokenized into integers** (Pre-process data)<br>- **Training data**: Pairs of tokens or sequence of tokens. Ex: User-item, sentence-sentence, **genre-description** | - **Shuffled JSON lines**<br>- **2 input channels**: 2 encoders and one comparator<br>- **Encoder types:** Average-pooled, CNN's, Bidirectional LSTM | - **Deep Learning**: dropout, epochs, learning_rate, batch_size, layers, optimizer, weight_decay, early_stopping<br>- **Enc1_network, Enc2_network**<br>  • Hcnn, blistm, pooled_embedding | - **Single Machine CPU**<br>- **Single Machine with Multi-GPU**<br>- **Inference:** ml.p2.2xlarge with INFERENCE_PREFERRED_MODE for encoding |
| Object Detection (WHERE?) | - **Single deep neural network, CNN with SSD**<br>- Identify **image within bounding boxes**<br>- **Labels to bounding boxes**<br>- **classes with confidence** scores<br>- new or pre-trained models from ImageNet | - RecordIO or image format<br>- JSON file for annotation data for each image | - **Input:** Image<br>- **Output:** all **instances of images** with **categories** and **confidence** scores<br>- **Transfer learning** | - Mini_batch_size, Learning_rate<br>- **Optimizer:** sgd, adam, rmsprop, adadelta | - **Training:** Multi-GPU or Multi-machine GPU (p2, p3)<br>- **Inference:** CPU or GPU for (C5, M5, P2, P3) |
| Image Classification (WHAT?) | - **RestNet CNN**<br>- **One or more labels to image**<br>- Tells **what objects** are | - MXNet RecordIO<br>- raw jpg or png images<br>- .lst files: associate image to index, class label, path to image<br>- **Reader**: Pipe mode to stream read images | - **Full training mode:** Initialize Network with random weights<br>- **Transfer learning:** Initialized weights, network fine-tuned with new training data<br>- **Default image size:** 3-channel 224x224 | - batch_size, Learning_rate<br>- **Optimizer**: weight decay, beta 1, beta 2, eps, gamma | - **Training:** Multi-GPU or Multi-machine GPU (p2, p3)<br>- **Inference:** CPU or GPU (C4, P2, P3) |
| Semantic Segmentation (Pixel-level) | - **MXNet Gluon and Gluon CV**<br>- **Pixel level image classification**<br>- Produce segmentation mask<br>- **Ex:** self-driving, medical | - JPEG and PNG annotations<br>- Training and validation<br>- Labels maps to describe<br>- **Augmented manifest image format** with Pipe mode | - **3 algorithms:** FCN, PSP, DeepLabV3<br>- **Backbones:** ResNet50, ResNet101 trained on ImageNet<br>- **Scratch or transfer learning** | - batch_size, Learning_rate, Optimizer, epochs<br>- Algorithm<br>- backbone | - **Training:** Only single machine GPU (P2 or P3)<br>- **Inference:** CPU (C5, M5) or GPU (P2, P3) |

| | | | | | |
|---|---|---|---|---|---|
| | imaging, robot sensing | - **JPG** images accepted for inference | | | |
| **Random Cut Forest** | - **Forest of trees (each tree partition of training data)**<br>- **Unsupervised Anomaly Detection**<br>- Unexpected dip or spike in time series data<br>- Assign anomaly score to each data point<br>- Developed by Amazon | - **RecordIO-protobuf or CSV**<br>- **Read**: File or Pipe<br>- **Optional test channel**: computing accuracy, precision, recall and F1 on labelled data | - **Expected change in complexity of tree when adding a point**<br>- Data is sampled randomly then trained, On Streaming data<br>- Works on Stream and Batch data | - **Num_trees:** Increasing reduces noise<br>- Num_samples_per_tree | - **Training:** CPU only (M4, C4, C5)<br>- **Inference:** CPU (ml.C5.xl) |
| **Neural Topic Model (Unsupervised for documents)** | - **"Neural Variational Inference"**<br>- **Unsupervised to classify or summarize documents**<br>- Organize docs into topics<br>- **Ex.** bike, car, train classify doc into transportation | - **RecordIO-protobuf or CSV**<br>- **4 channels:** train (mandatory), validation, test, auxiliary.<br>- **Words tokenized into integers:** Count of every word in doc in CSV vocabulary in auxiliary channel.<br>- **Read:** File or Pipe | - **Define how many topics**<br>- **Topics:** latent representation based of top ranking words | - **Lowering Mini_batch_size, Learning_rate:** reduce validation loss at expense of training time<br>- Num_topics | - **Training:** CPU or GPU(recommended)<br>- **Inference:** CPU |
| **Latent Dirichlet Allocation (Unsupervised topic model, not DL)** | - **Unsupervised Topic modelling with no DL**<br>- Group docs with shared subset of words<br>- **Ex**: Cluster customer based on purchase | - **RecordIO-protobuf or CSV**<br>- **2 channels:** Train, optional test<br>- Each doc with count of every word in CSV vocabulary<br>- **Read:** File or pipe(only with recordIO) | - **Define how many topics**<br>- **Optional Test channel:** scoring results<br>- **Similar to NTM but CPU based (cheaper and efficient)** | - Num_topics,<br>- **Alpha0** - Initial guess for concentration paramerer, small = sparse topics | - **Training:** Single-instance CPU |
| **KNN - K-Nearest-Neighbors (Supervised Classification)** | - **Supervised Classification or regression**<br>- **Classification:** K closest points to sample point and return most frequent label<br>- **Regression**: K closest points to sample point and return average | - **RecordIO-protobuf or CSV**<br>- **Train channel:** Data<br>- **Test channel:** Emit accuracy or MSE<br>- **Read:** File or Pipe | - **Data is first sampled**<br>- **Dimension reduction**<br>- **Build index for looking up neighbors**<br>- **Serialize model**<br>- **Query model for a given K** | - K = how many neighbors to look at<br>- Sample_size | - **Training:** CPU (C5) or GPU (P2)<br>- **Inference:** CPU (Low Latency), GPU (High throughput on large batches) |
| **K-Means Clustering (Unsupervised Segmentation, Grouping)** | - **Unsupervised clustering classification**<br>- Divide data into K groups<br>- member of groups similar to each other<br>- Web-scale K-means clustering | - **RecordIO-protobuf or CSV**<br>- **Train channel:** Data ShardedbyS3Key<br>- **Test channel:** FullyReplicated<br>- **Read:** File or Pipe | - Every observation mapped to n-dimensional space (n = number of features)<br>- **Optimize center of K clusters:** Determine initial, calculate from training data, reduce cluster | - **K:** Use "elbow method", optimize for tightness of clusters<br>- Mini_batch_size<br>- Extra_center_factor<br>- Init_method | - **Single GPU per instance (p*.xlarge)**<br>- **CPU(Recommended)** |
| **Principal Component Analysis (PCA) (Unsupervised)** | - **Unsupervised Dimensional reduction technique**<br>- minimize loss of info<br>- reduced dimensions called components: **1st comp with large variability,** second next largest<br>- **Ex:** Online surveys, census data | - **RecordIO-protobuf or CSV**<br>- **Read:** File or Pipe | - Covariance matrix created then SVD (singular value decomposition)<br>- **Regular mode algo:** sparse data and moderate observations and features<br>- **Randomized mode algo:** large observations and features and uses approximation algorithm | - Algorithm_mode<br>- Subtract_mean = Unbias data | - **GPU or CPU:** depends on input data |
| **t-Distributed Stochastic Neighbor Embedding (TSNE)** | - **non-linear dimensionality reduction algorithm** | | | | |
| **Factorization Machines (Supervised regression)** | - **Supervised Classification or regression of Sparse data**<br>- Used in context of **recommender system, click rate pattern**<br>- **Ex**: recommender system with less data from individual user<br>- Limited to pair-wise interactions. User --> item example | - **Must be RecordIO-protobuf with Float32** | - **Find factors we can use to predict classification**<br>- E**x: click or not click? Predicted rating? Users and items?** | - bias, factors, linear terms<br>- Uniform, normal or constant | - **CPU: recommended**<br>- **GPU** only for dense data |
| **Collaborative Filtering (Matrix Factorization)** | - **Intelligent recommender**<br>- filter out items that a user might like on the basis of reactions by similar users | | | | |
| **IP Insights (Unsupervised)** | - **Unsupervised neural network - IP address usage patterns**<br>- **identify suspicious behavior from given IP**: logins, accounts creating resources | - **CSV only (Entity,IP)**<br>- Directly feed Usernames, AccountIDs<br>- Training channel<br>- optional validation (computes AUC score) | - Neural network to **learn latent vector representation of entities and IP add**<br>- Hashed and embedded entities<br>- **Auto-generate negative samples** during training<br>- Random pairing entities and IP | - **Num_enity_vectors:** Hash size, 2x unique entity identifiers<br>- **Vector_dim:** Size of embedding vectors<br>- Epochs, learning_rate, batch_size.etc | - **CPU:** size depends on vector_dim and num_entity_vectors<br>- **GPU:** recommended, P3, can use Multiple GPU's |
| **Reinforcement Learning (DL)** | - **DL framework with Tensorflow and MXNet**<br>- **Agent that explores space**<br>- learn state changes in | | | - Specific **based on implementation**<br>- **Optimize using SM hyperp**arameters | - **DL so GPU**<br>- **Multi-core, multi-instance.**<br>- **Distributed training** |

| | | | | | |
|---|---|---|---|---|---|
| | different conditionsto predict subsequent behavior.<br>- **Ex:** Pac-Man, Supply chain, Autonomous vehicles, Industrial robotics | | | | **and/or environment rollout** |
| Q-Learning | - **Implementation of Reinforcement Learning**<br>- **Q = state/action** (S: environmental states, a: possible actions in states )<br>- Starts with zero then **penalty/reward based on learning space**<br>- Randomly explore choices<br>- **Use Q stored values to inform future choices** | | | | |