| Algorithm | Algo. Type | Input Format | INT/FLOAT | Processor Type | Instance | Multiprocessor in Single Machine | Multi Machine | Use Cases | Comments | HP |
|---|---|---|---|---|---|---|---|---|---|---|
| **Linear Learner** | SUPERVISED | * RecordIO **Wrapped** Protobuf / CSV<br>* Float32 Data only | FLOAT32 | CPU<br>GPU | Any | CPU<br>GPU | Only CPU<br>No GPU | * Regression and classification<br>* Classification: Binary or multiclass | * Need data to be normalized else algo may not converge<br>* Multiple models are trained in parallel | **balance_multiclass_weights**<br>learning_rate<br>mini_batch_Size<br>L1, L2 |
| **XGBoost** | SUPERVISED | * CSV or LibSVM (not AWS algo, but adapted, hence <span style="color:red">NO</span> RecordIO-protobuf) | - | CPU | M4 | - | No | * Regression and classification<br>* Classification: Binary or multiclass | * Output model as pickle<br>* Uses extreme boosting of trees<br>* Algo is memory bound, not much compute | * **subsample_trees** (less overfitting)<br>* eta (eq. to learning rate)<br>* alpha, gamma, lambda (conservative trees for higher values) |
| **Seq2Seq** | SUPERVISED | * RecordIO-Protobuf | INT | GPU | P3 | GPU | No | * Machine translation<br>* Text summarization<br>* Speech to text<br>* Any use case where input a sequence and output is a sequence | * Along with training data and validation data files, must provide **vocabulary files -- in case of text seq2seq**<br>* Start with tokenized text files, then convert to RecordIO-Protobuf<br>* Uses RNNs and CNNs internally | * batch_size<br>* optimizer<br>* learning_rate<br>* **num_layers_encoder**<br>* **num_layers_decoder**<br>* can optimize on: accuracy, BLEU score (mach. translation), perplexity |
| **DeepAR** | SUPERVISED | * JSON Lines<br>* GZIP<br>* Parquet<br><br>-- Each record to contain<br>- **Start**: starting TS<br>- **Target**: the TS values to learn/predict | - | CPU<br>GPU | C4<br>P3 | CPU<br>GPU | CPU<br>GPU | * Stock price prediction<br>* Sales and promotion effectiveness<br>* Any time oriented forecasting, single dimension | * Uses **RNNs**<br>* Can train several related timeseries, more series the better results, learns relationships b/w timeseries<br>* Start with CPU (C4.2xlarge, or higher), if necessary, move to GPU. Only large models need GPU | * **context_length** (number of time points back in time the model learns)<br>* epochs, batch_size, learning_rate, num_cells |
| **Blazing Text - Text Classification** | SUPERVISED | Augmented manifest text format --<br>"__label__1 this is a sentence with , punctuations also tokenized . that is space delimited . One sentence per line . label at the start" | - | CPU<br>GPU | size < 2GB: C5<br>size > 2GB: P2, P3 | Single GPU | No | * web search and information retrieval | * predict labels for sentence | * epochs<br>* learning_rate<br>* word_ngrams<br>* **vector_dim** |

| Algorithm | Algo. Type | Input Format | INT/FLOAT | Processor Type | Instance | Multiprocessor in Single Machine | Multi Machine | Use Cases | Comments | HP |
|---|---|---|---|---|---|---|---|---|---|---|
| **Blazing Text - Word2Vec** | UNSUPERVISED | **Word2Vec**<br>one sentence per line | - | CPU<br>GPU | P3 | CPU/GPU: CBOW & Skip Gram | GPU: Batch skip gram<br>CPU: No | * Preparing input for NLP use cases<br>* Vectorization of text for machine translation and sentiment analysis<br>* Semantic similarity of words | * Represents words as vectors<br>* Semantically similar words are represented by vectors close to each other<br>* Semantic -- of or relating to meaning in language<br><br>**MULTIPLE MODES:**<br>* *CBOW* - Continuous Bag of Words - Order of words DO NOT matter<br>* *Skip Gram* i.e. n-gram - order of words matter<br>* *Batch skip gram* - order of words matter | * mode: mandatory<br>* learning_rate<br>* window_size<br>* **vector_dim**<br>* negative_Samples |
| **Object2Vec** | | * Any object to be tokenized into integers<br>* Training data:<br>- pairs of tokens<br>- sequence of tokens | INT | CPU<br>GPU | M5, P2 | Single machine | No | * Collaborative recommendation system<br>* Multi-label document classification system<br>* Sentence Embeddings<br>* Learns relations or associations:<br>- sen to sen<br>- labels to seq (genre to description)<br>- product to product (recommendation)<br>- user to item (recommendation) | * CNNs and RNNs used<br>* Encoders used in input<br>- uses 2 encoders in parallel<br>- learns associations b/w encoders, using a comparator<br><br>**Encoder types:**<br>* Hierchical CNNs (hCNNs)<br>* bi-lstm<br>* pooled_embedding | dropout, early_stopping_epochs, learning_rate, batch_size, layers, act. func., optimizer, weight_decay |
| **Object Detection** | SUPERVISED | RecordIO (NOT Protobuf) or Images (JPEG or PNG)<br>+<br>With image manifest in JSON, one JSON per image that contains annotations | - | GPU | P2, P3 | Yes | Yes | * Detect objects in an image<br>* Object tracking | * Uses CNN with SSD<br>* Transfer learning/incremental learning supported<br>* Uses FLIP, RESCALE, JITTER internally to avoid overfitting<br>* CPUs can be used for inference, not for training | Standard CNN HPs like: learning_rate, batch size, optimizer etc. |
| **Image Classification** | SUPERVISED | * Pipe: Apache MxNET RecordIO (NOT Protobuf) - for interoperability with other DNN frameworks<br>* File Mode: Raw JPEG/PNG + *.LST files - associates image index, class label, path to image<br>-- To use images directly in Pipe mode use JSON based Augmented Manifest Format | - | GPU | P2, P3 | Yes | Yes | * classify images into multiple classes<br>* dog/cat/rat/tiger etc. | * **Full training**: ResNET CNN is used. N/W initialized with random weights<br>* **Transfer Learning/Pre-trained**: Image Net is used. Initialized with pre trained weights. Only Top FC layer is initialized with random weights.<br><br>* CPU can be used for inference, if not suitable, move to GPU | * batch_size<br>* learning_rate<br>* optimizer, B1, B2, eps, Gamma |

| Algorithm | Algo. Type | Input Format | INT/FLOAT | Processor Type | Instance | Multiprocessor in Single Machine | Multi Machine | Use Cases | Comments | HP |
|---|---|---|---|---|---|---|---|---|---|---|
| **Semantic Segmentation** | SUPERVISED | * Raw JPEG/PNG in file mode + annotations<br>* Add Augmented Manifest Format for Pipe Mode | - | GPU | P2, P3 | Yes | No | * Self driving cars<br>* Medical imaging and diagnostics<br>* Robot sensing<br>* Given a pixel - what object does it belong to ? | * Algo under hood: Gluon CV of MxNET = FC + Pyramid Scene Pairing + DeepLabV3<br>* Arch: **ResNet50/ResNet101** = "**Backbone**" selection in HP<br>* Trained on ImageNet data<br>* Incremental/Transfer learning allowed<br>* Inference can use CPU or GPU<br><br>Each of the three algorithms has two distinct components:<br>* The **backbone (or encoder)**—A network that produces reliable **activation maps** of features.<br>* The **decoder**—A network that constructs the **segmentation mask from the encoded activation maps**.<br><br>The segmentation output is represented as a **grayscale** image, called a **segmentation mask. A** segmentation mask is a grayscale image with the **same shape as the input image.** | epochs, learning_rate, batch size, algo, **backbone** |
| **Random Cut Forest** | UNSUPERVISED | * RecordIO-Protobuf<br>* CSV | - | CPU | M4,C4,C5 | - | No | * Anomaly detection<br>* Detect unexpected spikes in TS data<br>* Few people have tried using this for fraud detection | * Assigns anomaly score to each data point<br>* Uses forest of trees<br>* Looks at **expected change in complexity as a result of adding a point to a tree**<br>* Random sampling<br>* RCF is used in Kinesis Analytics in real time | num_trees,<br>**num_samples_per_tree** (= choose inversely proportional to ratio #anomalous/#normal in dataset) |
| **Neural Topic Modelling** | UNSUPERVISED | * RecordIO-Protobuf<br>* CSV<br><br>- Words must be tokenized to integers<br>- aux channel for vocab | INT | GPU | P2, P3 | - | | * Organize docs into topics<br>* Summarize docs based on topics | * Algo: **Neural Variational Inference**<br>* Define how many topics to group docs into<br>* Used only on text<br>* CPU / GPU for inference | **num_topics**<br>mini_batch_size<br>learning_rate<br>**variation_loss** (at expence of learing time) |
| **LDA (Latent Dirichlet Allocation)** | UNSUPERVISED | * RecordIO-Protobuf (Pipe Mode)<br>* CSV<br><br>- Words must be tokenized to integers<br>- aux channel for vocab | - | CPU | M4 | No | No | * Cluster customers based on purchases<br>* Harmonic analysis in music | * Algo: LDA - Open source availability, **not DNN**<br>* Can process more than text, like harmonic music analysis<br>* Single inst. CPU | **num_topics**<br>**alpha0** = small values - sparse topic mixtures, >1 uniform topic mixture |

| Algorithm | Algo. Type | Input Format | INT/FLOAT | Processor Type | Instance | Multiprocessor in Single Machine | Multi Machine | Use Cases | Comments | HP |
|---|---|---|---|---|---|---|---|---|---|---|
| **kNN (k Nearest Neighbors)** | SUPERVISED | * RecordIO-protobuf<br>* CSV<br>-- File or pipe mode both<br>- first column has label | - | CPU<br>GPU | - | - | - | * Classification and regression | * Sagemaker automates 3 steps:<br>- Sample data (can't use for huge data)<br>- Dim reduction (sign or nfjlt methods)<br>- Build index for looking up neighbours | **k**<br>**sample_size** |
| **K-Means** | UNSUPERVISED | * RecordIO-protobuf<br>* CSV<br>-- File or pipe mode both | - | CPU (recommended)<br>GPU | M4, M5, C4, C5 | - | - | * Cluster data - unsupervised<br>* Find groups of data points based on similarity | * Webscale K-Means in Sagemaker<br>* Similarity measured by euclidean distance<br>* Works to optimize the centers of eack of the k-clusters<br>* **Algorithm**:<br>1) Determine init. cluster centers = 2 ways: **k-means++** (tries to make initial clusters far apart) OR **random**<br>2) Iterate over data and calculate cluster center<br>3) Reduce from K to k - using Lloyd's method or k-means++<br><br>K comes from "extra_cluster_centers" which improves accuracy, but later reduced to k.<br>K = k * x | * **K**<br>* mini_batch_size<br>* **extra_center_factor (x)**<br>* **init_method** (k-means++ OR random) |
| **PCA - Principal Component Analysis** | UNSUPERVISED | * RecordIO-protobuf<br>* CSV<br>-- File or pipe mode both | - | CPU<br>GPU | - | - | - | * Dimensionality Reduction<br>* Removes Curse of Dimensionality | * Reduced Dimensions are called components<br>* 1st component - largest possible variaility, next 2nd component, so on ..<br>* Used Singular Value Decomposition (SVD)<br>* **Two Modes**:<br>- **Regular**: Sparse data. modelate #features, #rows<br>- **Randomized**: Dense data. #large data, #large features, uses approximation algos | * **algorithm_mode** (regular, random)<br>* **subtract_mean**: unbiases data |
| **Factorization Machines** | SUPERVISED | * RecordIO-Protobuf | FLOAT32 | CPU (recommended)<br>GPU | - | - | - | * Regression, Classification, recommendation - all in one general purpose algo for sparse data<br>* Click prediction<br>* Item recommendation | * Limited to pairwise interaction - 2nd order<br>e.g. user to item interactions<br>* **CSV not practical** hence not supported,a s data is sparse<br>* **GPU not recommented** as data is sparse, GPU works better on dense data | * **Initialization methods** for bias, factors and linear terms<br>- methods: uniform, normal or const<br>- can tune properties of each method |

| Algorithm | Algo. Type | Input Format | INT/FLOAT | Processor Type | Instance | Multiprocessor in Single Machine | Multi Machine | Use Cases | Comments | HP |
|---|---|---|---|---|---|---|---|---|---|---|
| IP Insights | UNSUPERVISED | * CSV only for training<br>* Inference: JSON lines, CSV, JSON | - | CPU GPU (recommended) | - | Multi GPU | - | * Identify suspicious IP addresses in context of security<br>* Logins from anomalous IPs<br>* Identify accounts creating resources from anamolous IPs | * **Only IPv4 supported**<br>* Uses **NN** to learn latent vector rep. of entities and IP addresses<br>* entities are hashed and embedded - large hash size<br>* Automatically generates anomalous data by randomly pairing entities and IPs - as data will be highly imbalanced | * **num_entity_vectors** (hash size, set to twice the unique entity identifiers)<br>* **vector_dim** (size of embedding vectors, scales model size)<br>* Others: epoch, batch_size, leraning_rate, etc. |
| Reinforcement Learning | REINFORCEMENT LEARNING | * Nothing specific to Sagemaker | - | GPU | GPU | Yes | Yes - Multi Instance GPU recommended | * Games<br>* Supply chain management<br>* HVAC Systems<br>* Industrial robotics<br>* Dialog systems<br>* Autonomous vehicles | * Supports Intel coach, Ray RLLib<br>* Tensorflow, MxNET<br>* Custom, commercial and opensource environments supported - Matlab simulink, energy plus, robo school, pybullet, Amazon Sumerian, AWS Robomaker | * Depends on framework and algo used, nothing tied to Sagemaker |

| Mandatory FLOAT32 | Mandatory INT32 | CPU ONLY | GPU Only | Incremental Training Available |
|---|---|---|---|---|
| Linear Learner | Seq2Seq | XGBoost | Seq2Seq | Image Classification |
| Factorization Machines | Object2Vec | RCF | Image Classification | Semantic Segmentation |
| | NTM | LDA | Semantic Segmentation | Object Detection |
| | | | Object Detection | |
| | | | NTM | |
| | | | RL | |