



USE OF SPARK MLLIB FOR PREDICTING OFFLINING OF DIGITAL MEDIA

Christopher Burdorf
NBCUniversal



OVERVIEW

- λ Problem Definition
- λ Cluster Configurations
- λ Parameters
- λ MLLib libraries utilized
- λ Results
- λ Conclusions



Problem Definition

- λ Digital media files distributed internationally 24/7 over internet to cable TV channels in APAC, Europe, and Asia.
- λ TV Shows, Commercials
- λ Fast Isilon storage fills up frequently, thus offlining is necessary to create space for new files
- λ Multiple parameters are used to pick candidates
- λ Previous system works well, but is slow and has large overhead



Cluster Configurations

λ Development

- 1 Linux System with 16 cores and 16 Gigs RAM
- 3 Mac OS Systems each with 8 cores and 16 Gigs RAM
- Mesos cluster manager

λ Production

- 3 Linux systems with 32 cores and 64 Gigs of RAM
- Mesos master and slaves running in Docker containers
-



Features/Parameters

- λ File size
- λ File age
- λ Days since last airing
- λ Days until next airing
- λ Immediately remove files that have not been scheduled for more than 3 weeks (was 2 weeks to start with).



K-Means Clustering

$$J = \sum_{j=1}^k \sum_{i=1}^n \left\| x_i^{(j)} - c_j \right\|^2$$

where $x_i^{(j)} - c_j$ is a chosen distance measure between a data point and the cluster centre , is an indicator of the distance of the n data points from their respective cluster centers.



K means clustering results

- λ Attempted with multiple centroid points
- λ No meaningful clusters for any sets of data attempted
- λ Results appeared to be tied to the nature of the data
- λ Multiple parameters created difficulties



Naive Bayes classifier

- λ From Bayes theorem: $p(C_k|x) = (p(C_k)p(x|C_k))/p(x)$
- λ posterior = (prior X likelihood)/evidence
- λ Thus, the attempt was to classify media files for removal or retention on the file system based on the different parameter values and their expected results
- λ Had difficulties training for many factors
 - Adding new factors to training set altered results for previously trained classifier and required retraining the entire classification system



Support Vector Machines

To obtain the optimal hyperplane, one solves the following convex quadratic optimization problem with respect to weight vector w and bias b :

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i, \quad (1)$$

subject to the constraints:

$$y_i(\langle w, \phi(x_i) \rangle + b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0 \quad \text{for } i = 1, \dots, n \quad (2)$$

The regularization parameter C controls the trade-off between maximizing the margin $1/\|w\|$ and minimizing the sum of slack variables of the training examples

$$\epsilon_i = \max(0, 1 - y_i(\langle w, \phi(x_i) \rangle + b)) \text{ for } i = 1, \dots, n. \quad (3)$$

The training example x_i is correctly classified if $0 \leq \epsilon_i < 1$ and is misclassified when $\epsilon_i \geq 1$.





SVM Results

- λ Found to be most robust method
- λ Compensated well for new added features
- λ Built training set from production data and auto-generated data that fit the criteria
- λ Spark Mllib optimization allows us to build a predictive system in just over an hour
 - Run twice daily due to constant changes in schedules and online media.



Production environment

- λ Spark Mllib SVM is trained and predictions are generated for each media file online twice daily. Predictions are stored in HBase
- λ Media manager daemon regularly scans the file system and if available space requires purging, it will select files based on predictions stored in Hbase
- λ Backup system checks to see if schedule has changed for a given media file selected to be purged, because Spark SVM predictions are only generated twice daily



Production issues

- λ System was run in test mode for 6 weeks without a problem
- λ Once switched to production, issues arose
 - 2 weeks was too short a time to immediately offline files
 - λ Switched to 3 weeks
 - Docker containers filled up with Spark temp files and crashed over the weekend (ouch!)
 - λ Solved with cron job periodically removing them.
 - Web service access to Hbase locked up.
 - λ Solved by having thread timeout on web services call.



Conclusions

- λ System has been running in production mode for some time now.
 - Fine tuning appears to be complete
- λ Spark SVM has performed well
 - Fast and robust
- λ Good application for machine learning though critical aspect of the system increased task complexity.