

Machine Learning Tutorial

Wrap-up

Jake Hofman

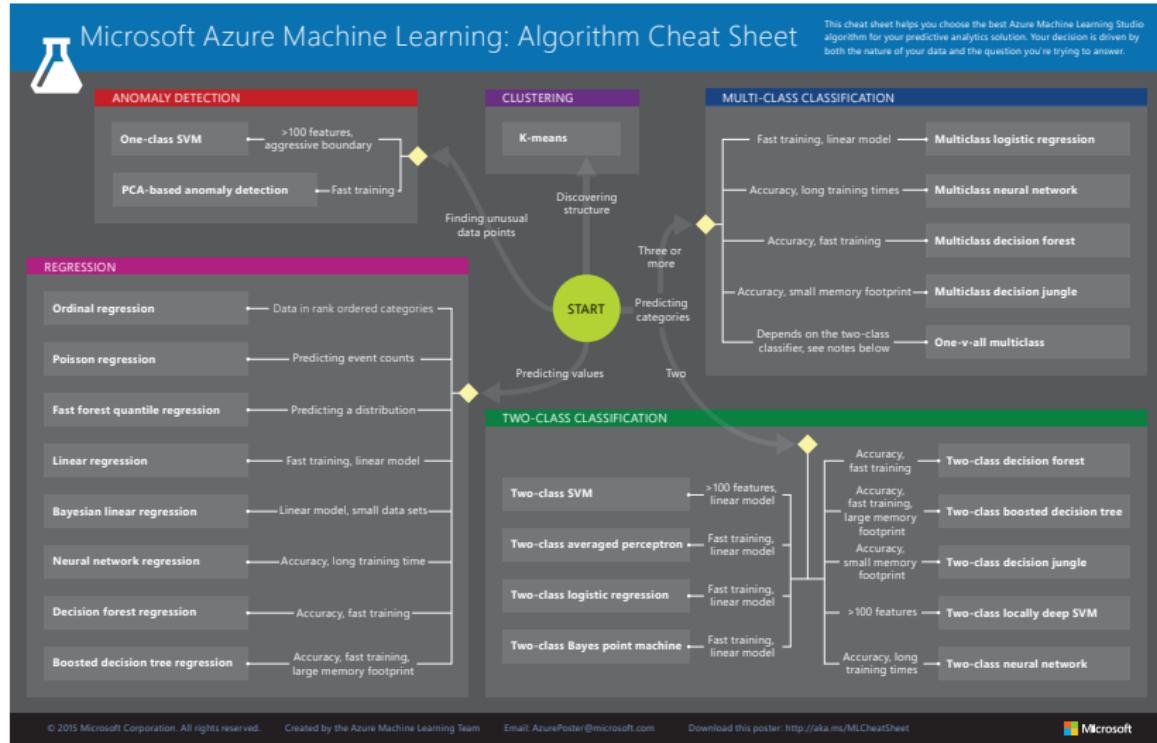
Microsoft Research

March 16, 2018

Books



Cheatsheets



© 2015 Microsoft Corporation. All rights reserved.

Created by the Azure Machine Learning Team

Email: AzurePoster@microsoft.com

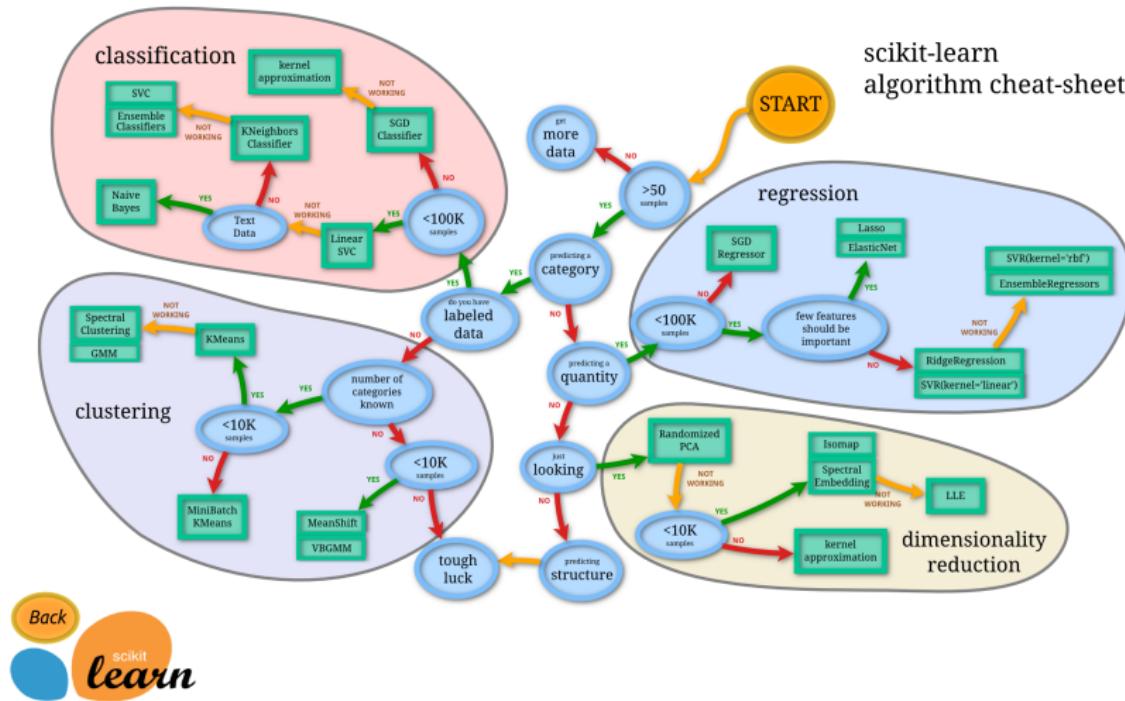
Download this poster: <http://aka.ms/MLCheatSheet>



<http://bit.ly/howtochooseml>



Cheatsheets



<http://bit.ly/scikitcheat>

Cheatsheets

CHEAT SHEET: ALGORITHMS FOR SUPERVISED- AND UNSUPERVISED LEARNING 1

Algorithm	Description	Model	Objective	Training	Decomposition	Decomposition	Update	Online algorithm
Labeled neighbor	The label of a new point x is classified with the most frequent label \hat{y} of the k nearest training instances.	$\hat{y} = \text{argmax}_{\hat{y}} \sum_{x_i \in \text{NN}(x)} \delta(y_i, \hat{y})$	No optimization needed.	For one-validation to learn the aggregate \hat{y} , otherwise no training, classification based on existing points.	\hat{y} acts as an regularizer the classifier, as $k \rightarrow N$ the boundary becomes smoother.	$O(N^2)$ space complexity, slow all training instances and all other features need to be kept in memory.	Only needs non-linear boundaries.	To be added.
Naive Bayes	Learn $p(x_i y)$ by dividing $p(x_i y)$ and $p(x_i)$ using Bayes rule. Then the class conditional probability is proportional to the joint probability of all others, e.g. $p(y x) = \frac{\prod_i p(x_i y) \cdot p(y)}{\sum_y \prod_i p(x_i y) \cdot \log(p(y))}$.	$p(y x) = \text{argmax}_y p(y x)$ $= \text{argmax}_y p(y) \cdot p(x y)$ $= \text{argmax}_y \prod_i p(x_i y) \cdot p(y)$ $= \text{argmax}_y \frac{\prod_i p(x_i y) \cdot \log(p(y))}{\sum_y \prod_i p(x_i y) \cdot \log(p(y))}$	No optimization needed.	Maximum Likelihood $p(\theta x) = \prod_{i=1}^{n_{\text{train}}} \log(p(x_i \theta))$ posterior $= p(\theta x) = \prod_{i=1}^{n_{\text{train}}} p(x_i \theta) \cdot p(\theta)$	Use a Dirichlet prior on the parameters to obtain MAP estimate.	$p(\theta x) = \text{argmax}_{\theta} \psi(\theta) + \sum_{i=1}^{n_{\text{train}}} \log(p(x_i \theta)) - \theta_0 \cdot \alpha - \theta_1 \cdot \beta - \theta_2 \cdot \gamma$	$O(N^2)$ space complexity, stock training instance must be visited and each of its feature research.	Can only learn linear boundaries for multinomial/multinomial distributions. With Gaussian features, quadratic space complexity with unmodelled distributions.
Log-Linear	Estimate $p(x y)$ directly by learning a maximum entropy distribution over the conditional entropy distribution.	$p(y x) = \text{argmax}_y p(y x)$ $= \text{argmax}_y \prod_{i=1}^{n_{\text{train}}} \log(p(x_i y))$... where $p(y x) = \frac{1}{Z} \sum_{\theta} \exp \left(\sum_{i=1}^{n_{\text{train}}} \theta_i \cdot h_i(x_i, y) \right)$ $Z = \prod_{i=1}^{n_{\text{train}}} \exp \left(\sum_{j=1}^{n_{\text{classes}}} \theta_j \cdot h_j(x_i, y) \right)$	Minimize the negative log-likelihood: $\mathcal{L}_{\text{log-lik}}(\theta, Z) = -\sum_{i=1}^{n_{\text{train}}} \log(p(x_i y))$ $\mathcal{L}_{\text{log-lik}}(\theta, Z) = -\sum_{i=1}^{n_{\text{train}}} \left(\log(p(x_i y)) - \sum_{j=1}^{n_{\text{classes}}} \theta_j \cdot h_j(x_i, y) \right)$... where θ_j is the set of weight j in test sample y ... where $h_j(x_i, y)$ are the indicator counts ... where $Z = \prod_{i=1}^{n_{\text{train}}} \exp \left(\sum_{j=1}^{n_{\text{classes}}} \theta_j \cdot h_j(x_i, y) \right)$	Gradient descent or gradient ascent of summing objective: $\theta_j = \theta_j - \eta \cdot \frac{\partial \mathcal{L}_{\text{log-lik}}}{\partial \theta_j}$... where η is the step parameter.	Possible large values for the λ parameter, try different values for λ and α to find a good fit.	$\mathcal{L}_{\text{log-lik}}(\theta, Z) = -\sum_{i=1}^{n_{\text{train}}} \log(p(x_i y))$ $\mathcal{L}_{\text{log-lik}}(\theta, Z) = -\sum_{i=1}^{n_{\text{train}}} \left(\log(p(x_i y)) - \sum_{j=1}^{n_{\text{classes}}} \theta_j \cdot h_j(x_i, y) \right)$... where θ_j is the set of weight j in test sample y ... where $h_j(x_i, y)$ are the indicator counts For each class c_j : $\sum_{i=1}^{n_{\text{train}}} \delta_{ij}(y_i) = \sum_{i=1}^{n_{\text{train}}} \delta_{ij}(p(x_i y))$	$O(N^2)$ space complexity, stock training instance must be visited and each of its feature research.	Approximate the class conditional distribution in terms of a kernel $K(x, x')$ and a weight vector w ($\psi(x) = 1 + w^T K(x)$). By the Representer Theorem $\psi(x) = \sum_{i=1}^{n_{\text{train}}} \alpha_i K(x, x_i) + b$. $p(y x) = \text{argmax}_y \left(\sum_{i=1}^{n_{\text{train}}} \alpha_i K(x, x_i) + b \right)$
Perceptron	Directly estimate the linear function $g(x)$ by iteratively updating the weights until no misclassified training instances remain.	$g(x) = \text{argmax}_y p(x y)$... where $p(x y) = \frac{1}{1 + e^{-y \cdot g(x)}}$... where $g(x) = \sum_{i=1}^{n_{\text{features}}} w_i \cdot x_i$	Try to minimize the linear function, the number of correctly classified input vectors.	Initial linear function $r_0 = 0$ and update the weight vector w after each iteration:	The Final Perceptron: the parameter w and the final weight vector's weight vector. Then $(w)^T \cdot \text{sign}(w^T x)$	$(w)^T \cdot \text{sign}(w^T x)$	Use a kernel $K(x, x')$ and b weight per training instance.	The perceptron is an online algorithm per update.
Support vector machines	A maximum margin classifier finds the separating hyperplane with the largest margin m to the closest data points.	Primal $\min_{w, b} \frac{1}{2} \ w\ ^2$ s.t. $\sum_{i=1}^{n_{\text{train}}} \delta_{ij}(y_i \cdot w + b) \geq 1 - \epsilon_i$ Dual $D(w) = \sum_{i=1}^{n_{\text{train}}} \sum_{j=1}^{n_{\text{train}}} \delta_{ij} \cdot y_i \cdot y_j \cdot w_i \cdot w_j$... where $\delta_{ij} = \sum_{k=1}^{n_{\text{classes}}} \delta_{ijk}$	\bullet Quadratic Programming (QP) \bullet SMO: Sequential Minimal Optimization (Bottou)	The soft margin SVM predictor is obtained by the dual solution and linear combination of the support vectors.	Primal $\min_{w, b} \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^{n_{\text{train}}} \delta_{ij}$ s.t. $\sum_{i=1}^{n_{\text{train}}} \delta_{ij} \cdot y_i \cdot w + b \geq 1 - \epsilon_i$ Dual $D(w) = \sum_{i=1}^{n_{\text{train}}} \sum_{j=1}^{n_{\text{train}}} \sum_{k=1}^{n_{\text{train}}} \sum_{l=1}^{n_{\text{train}}} \delta_{ijkl} \cdot y_i \cdot y_j \cdot y_k \cdot y_l \cdot w_i \cdot w_j \cdot w_k \cdot w_l$ s.t. $0 \leq w_i \leq C, \sum_{i=1}^{n_{\text{train}}} \delta_{ij} \cdot w_i = 0$	\bullet QP (QP3)	Use a non-linear kernel $K(x, x')$ $p(x) = \sum_{i=1}^{n_{\text{train}}} \alpha_i K(x, x_i) + b_0$	Online SVM, see for example The Multi-class Linear and Kernel SVM and Perceptron
Linear	A hard margin, geometric clustering algorithm that finds a line that splits the data points in two classes A and B .	Separation: $\epsilon_{\text{sep}} = \frac{1}{2} \cdot \frac{ w }{\sqrt{1 + w^T w}}$ Minimization: $\frac{w^T w}{2} + C \sum_{i=1}^{n_{\text{train}}} \delta_{ij}(y_i \cdot w + b)$... i.e. minimize the distance from each cluster center to each of the points.	Separation: $\epsilon_{\text{sep}} = \frac{1}{2} \cdot \frac{ w }{\sqrt{1 + w^T w}}$ minimal for δ_{sep} Minimization: $\frac{w^T w}{2} + C \sum_{i=1}^{n_{\text{train}}} \delta_{ij}(y_i \cdot w + b)$... where δ_{sep} is the control of cluster k .	Only hard margin assignment to clusters.	To be added.	For non-linearly separable data, no kernel functions as required.	Support vector machines update the weights after processing one point at a time.	
Mixture of Gaussians	A probabilistic clustering algorithm that finds a set of Gaussian and their data points to represent the observed data from a particular Gaussian.	$\text{Approximation to clusters by specifying Gaussian components:}$ $p(x w, \mu, \Sigma) = p(\mu_1 w_1) \cdot p(\mu_2 w_2) \cdots$... with $w = (\text{Mixture}(\mu), \Sigma)$, and $\text{Mixture}(\mu) = \sum_{i=1}^{n_{\text{mixtures}}} w_i$ is the total weight of the mixture and w_i is the weight of the i -th Gaussian component.	Log-likelihood: $\mathcal{L}(w, \mu, \Sigma X) = \log(p(X w, \mu, \Sigma))$ $= \sum_{i=1}^{n_{\text{train}}} \log \left(\sum_{j=1}^{n_{\text{mixtures}}} w_j \cdot \text{prob}(x_i \mu_j, \Sigma_j) \right)$	Maximization: $\frac{\partial \mathcal{L}}{\partial w} = \frac{1}{2n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \sum_{j=1}^{n_{\text{mixtures}}} \frac{\partial \log(p(x_i w, \mu, \Sigma))}{\partial w_j}$ $w_j = \frac{\sum_{i=1}^{n_{\text{train}}} \delta_{ij}}{\sum_{i=1}^{n_{\text{train}}} \delta_{ij}}$	The mixture of Gaussians assigns probabilities for each cluster to each data point, and so each is capable of capturing ambiguities in the data set.	To be added.	Not applicable.	Online Gaussian Mixture Models: A good start to: A View of the EM Algorithm that Justifies Mixture Models (Dempster, Laird & Rubin, 1977).

*Created by Eferman Hofman, 07/2011, for semi-practitioners reasons while studying for a Machine Learning exam. Last updated May 1, 2012.

soulmachine@gmail.com

Machine Learning Cheat Sheet

Classical equations, diagrams and tricks in machine learning

February 12, 2015

<http://bit.ly/mlcheatbook>

Applied ML Course

Data-driven modeling

Spring 2012, Department of Applied Mathematics, Columbia University

[« Home](#)

[Pages](#)

[Instructor](#)

[Overview](#)

[Recent Posts](#)

[Syllabus](#)

[RSS Feeds](#)

[All posts](#) 

[All comments](#) 

[Search](#)

[Find](#)

OVERVIEW

DATA-DRIVEN MODELING

Spring 2012

Department: Applied Mathematics, Columbia University

Instructor: [Jake Hofman](#)

Course number: E4990

Time: Mondays, 4:10-6:40pm

Location: 627 Seeley W. Mudd Building

Overview:

This course is an introduction to applied problems in statistics and machine learning. Lectures will cover the theory behind simple but effective methods for supervised and unsupervised learning as well as tools and techniques for acquiring, cleaning, and utilizing data to solve real-world problems. Emphasis will be on formulating real-world modeling and prediction tasks as optimization problems and comparing methods in terms of practical efficacy and scalability. Students will gain direct experience in acquiring data from online sources and will develop the necessary computing skills to address problems such as spam filtering and recommendation systems. The course will also feature guest lectures from prominent local practitioners in academia and industry highlighting how these skills are used in both research and business settings.

<http://jakehofman.com/ddm>



Modeling Social Data

Spring 2017

Department of Applied Physics and Applied Mathematics

Columbia University

Instructor: Jake Hofman

Course Numbers: [Applied Mathematics E4990](#)

Time: Fridays, 10:10am-12:40pm

Location: [633 Seeley W. Mudd Building](#)

[310 Fayerweather Hall](#)

This class focuses on data-driven models for social data—data that capture how people behave and interact with each other or with online platforms. The course will focus on the challenges that arise when working with large-scale observational data. We will present data science and data engineering methods needed for analyzing such real-world data at scale, focusing on learning models which balance predictive power and interpretability. In addition to core computational and statistical concepts, the course will also address practical issues around collecting, manipulating, and analyzing data with APIs, Unix tools, and statistical programming libraries.

<http://modelingsocialdata.org>

MOOCs

The screenshot shows the DataCamp website with a teal header. The header features the DataCamp logo (a white icon of a person inside a shield-like shape) and the word "DataCamp". To the right of the logo are four navigation items: "Learn ▾", "Pricing", "For Business", and "About ▾". Below the header, the text "SKILL TRACK" is displayed in a yellow font. The main title "Machine Learning with R" is prominently shown in large, bold, white text. A descriptive paragraph below the title explains the track's content: "Learn the basics of prediction using machine learning. This track covers predicting categorical and numeric responses via classification and regression, and discovering the hidden structure of datasets (unsupervised learning). Learn how to process data for modeling, how to train your models, how to visualize your models and assess their performance, and how to tune their parameters for better performance." At the bottom left of the main section is a yellow button with the word "Join" in white.

<https://www.datacamp.com/tracks/machine-learning>

MOOCs

The screenshot shows the DataCamp website with a teal header. The header features the DataCamp logo (a stylized icon inside a hexagon followed by the word "DataCamp"). Navigation links include "Learn", "Pricing", "For Business", and "About". Below the header, the text "SKILL TRACK" is in small yellow capital letters, followed by a large white title "Machine Learning with Python". A descriptive paragraph below the title states: "Machine learning is changing the world and if you want to be a part of the ML revolution, this is a great place to start! In this track, you'll learn the fundamental concepts in Machine Learning." At the bottom left of this section is a yellow button with the word "Join" in black text.

[http://www.datacamp.com/tracks/
machine-learning-with-python](http://www.datacamp.com/tracks/machine-learning-with-python)

MOOCs

The screenshot shows the Coursera website interface. At the top, there is a navigation bar with the Coursera logo, a 'Catalog' button, a search bar, and links for 'For Enterprise' and 'Log In'. Below the navigation bar, the main content area displays the 'Machine Learning' course page. The page title 'Machine Learning' is prominently displayed in large white text on a dark background. Above the title, there are breadcrumb links: 'Home > Data Science > Machine Learning'. To the left of the main content, there is a sidebar with links for 'Overview', 'Syllabus', 'FAQs', 'Creators', and 'Ratings and Reviews'. A large blue button labeled 'Enroll' with the text 'Starts Mar 19' is visible. Below the sidebar, there is a link 'Apply for Financial Aid'. The main content area contains a detailed description of the course, mentioning machine learning as the science of getting computers to act without being explicitly programmed. It highlights recent developments like self-driving cars, speech recognition, and improved genome understanding. A 'More' link is present at the bottom of this section. Below this, it says 'Created by: Stanford University' with a small Stanford University logo. Further down, there is a circular profile picture of Andrew Ng and text stating 'Taught by: Andrew Ng, Co-founder, Coursera; Adjunct Professor, Stanford University; formerly head of Baidu AI Group/Google Brain'.

<https://www.coursera.org/learn/machine-learning>

Course++

The screenshot shows a website for the CILVR Lab @ NYU. The header features a purple navigation bar with a molecular structure icon, the text "CILVR Lab @ NYU", and links for "People" and "Courses". On the left, there's a sidebar with links for Home, News, Events, Publications, People, Research, Software, Data, Courses, Sponsors, Contact, and groups for Fergus, LeCun, and Sontag. A "internal" link is also present. The main content area has a title "Big Data, Large Scale Machine Learning" and a "Table of Contents" sidebar with links to Course Information, News, Course Material, Prerequisites, Syllabus, and Evaluation. The "Course Information" section contains details about the term (Spring 2013), class time (Tuesdays 5:00 to 6:50 pm), classroom location (Warren Weaver Hall Room 109), instructors (John Langford and Yann LeCun), teaching assistant (Xiang Zhang, xiang.zhang AT nyu.edu), and a discussion group on Piazza. The "News" section lists various announcements and changes, such as assignment releases, classroom changes, and lecture topics.

<http://bit.ly/nyubigdata>

Course++

Machine Learning the Future

This is class 1pm-2:15pm Mondays in Spring 2017 at Cornell Tech and [Cornell](#) about the frontier of machine learning. Each class is lecture and discussion around a chosen topic. I would like each student to leave with an understanding of the topic deep enough to think about the next step.

The class will cover topics broadly related to reinforcement learning including contextual bandits, imitation learning, and exploration as well as other select topics drawn from online learning, optimization, boosting, parallelization, logarithmic prediction, active learning, representation. Logistically, we'll use [zoom](#) in addition to a video link.

January 30: cancelled (My voice is not working)

February 6: ML the Future slides, Generalization slides with source, background reading, with a Youtube recording and a backup recording.

February 13: cancelled (travel)

February 20: No class (February break)

February 27: Online Linear Learning slides with source Youtube recording and a backup recording papers: Importance Invariant, Adagrad, Normalized Dataset: RCV1 CCAT-or-not

March 6: Contextual Bandit Eval and Optimization slides with source, A Youtube recording and a backup recording papers: Double robust policies Dataset: RCV1 CCAT-or-not in contextual bandit format (uniform exploration)

March 13: Contextual Bandit Exploration slides with source, A Youtube recording and a backup recording, Papers Epoch Greedy, Cover, Bootstrap Dataset: RCV1 CCAT-or-not in multiclass format

March 20: Nonstationary evaluation slides with source Decision Service slides Contextual Decision Process slides Youtube recording and a backup recording Nonstationary evaluation, Decision Service, and Bellman Rank papers Decision Service site

March 27: Joint prediction as learning to search, Youtube recording and a backup recording DAgger paper AggreVaTe paper

April 3: no class(Spring break)

April 10: Joint prediction results, analysis, programming slides with source, A Youtube recording and a backup recording, LOLS paper Credit Assignment Compiler paper

April 17: Cancelled (healing up)

April 24: Logarithmic time prediction with source Youtube recording and a backup recording covering Error-Correcting Tournaments, Label embedding trees Dynamic Trees I, and Dynamic Trees II. Not covered: Static contextual bandits, Dynamic log time class probability Efficient Multilabel

May 1: Active Learning with source a Youtube recording and a backup recording Iwai papers search papers Cost sensitive active learning

May 8: Parallel Learning with source a Youtube recording and a backup recording alfricredu paper disbelief paper 1-bit SGD Ring alfricredu

<http://hunch.net/~mltf>

Thanks. Questions?

jmh@microsoft.com
@jakehofman