

Interactive Machine Learning for Discovering Patterns in Spectral Data and Images. D. A. Oyen¹ and N. L. Lanza¹, ¹Los Alamos National Laboratory, Los Alamos, NM 87545, doyen@lanl.gov.

Introduction: Data analysis methods like machine learning are advancing the ability to process large amounts of data quickly, even for data sets that do not fit standard statistical models. While interesting science questions cannot be answered completely through automated data processing, machine learning can improve the science return of remote sensors by increasing the speed at which scientists explore answers to their questions. Interactive machine learning balances the strengths of machine learning to perform repetitive pattern recognition tasks, while empowering scientists to explore interesting patterns in large sets of data [1].

Spectrometers are increasingly used in remote sensing, yet spectral data can be difficult to analyze due to its high-dimensionality and non-linear mapping to interpretable quantities. As part of the Mars Science Laboratory rover operations, ChemCam's Laser-Induced Breakdown Spectroscopy (LIBS) instrument collects fine-scale atomic spectra from targets up to 7m away [2]. Given the high number of ChemCam observations to date (>300,000) and the high dimensionality of LIBS spectra (~6000 channels), advanced analysis methods are needed. We are developing an interactive machine learning algorithm for discovering surface compositional features on rocks [3] in ChemCam targets. Using the insight that the precision of element abundance is more reliable than accuracy [4], we bypass the quantification of elements, and look directly for patterns of chemical gradients [5, 6]. By starting with a summary of depth trends, a variety of patterns can be discovered before narrowing in on a specific signature through interactive machine learning.

Machine Learning: In mathematical terms, machine learning is a form of function approximation. Given data X , with label Y , the goal is to find a function f such that $Y = f(X, \theta)$. X can be a vector, a spatial image, a timeseries, a spectral response, or it can contain even more complex structure. Y is typically a categorical value, but could be a scalar, vector or even a structured output. The functional form of f is typically fixed, such as a support vector machine or

neural network. The machine learning algorithm finds the optimal values of the parameters, θ , of the function f to make $Y = f(X, \theta)$. Interactive machine learning allows a person to modify X , Y or adjust constraints on θ and quickly get updated machine learning results.

When machine learning is used for pattern discovery, we have data X but do not have labels Y . Therefore, *unsupervised* machine learning takes the form of probability density estimation, such that $X \sim P(\theta)$, for a fixed distribution family P and learned parameters θ . The structure of the probability distribution is typically the most interesting aspect because it reveals interesting patterns about the data. Some examples include clustering which assumes that P is a distribution with multiple modes (or centers of clusters); and probabilistic graphical models which assume that P is a multi-variety joint distribution that can be factored compactly indicating direct dependencies.

Gaussian Graphical Models: Probabilistic graphical models [7], and specifically, Gaussian graphical models (GGM) [8], are unsupervised learning models that assume that each data sample $X = (x_1, x_2, \dots, x_p)$ is a p -dimensional vector generated by a multivariate joint distribution. Furthermore, the probability distribution can be factored into a compact representation with just a few direct dependencies. The compact representation assumption is a statistical necessity for the robust estimation of a high-dimensional distribution from finite data; and it reveals interesting structure about the dependencies among variables. Graphical models are used extensively in fields such as biology to identify gene interaction networks and neuroscience to infer functional pathways of neural activity.

To analyze the depth trend of a rock target at a location, we estimate *partial correlations* among spectra using the GGM algorithm. A partial correlation between shot A and shot B is the residual correlation after accounting for all other shots. Thus, a partial correlation is an estimate of a direct dependency. If the partial correlation between A and B is 0 then A and B are conditionally independent. A GGM is estimated

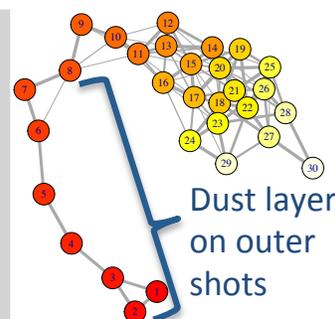
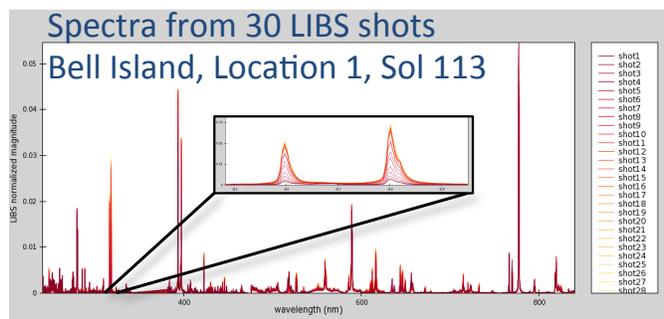


Figure 1 Interactive machine learning takes spectral data from a several LIBS shots (**left**) and learns a GGM (**right**) indicating geochemical trends in ChemCam targets.

from a data matrix X , where each column X_j is a shot j with spectral values X_{ij} for i in $\{1, \dots, n\}$ wavelengths. The sample covariance matrix, Σ , is calculated from X , then the best sparse approximation, Θ , to the partial correlation matrix for a given sparsity constraint, λ , is estimated. The number of non-zero partial correlations is controlled by the value of λ , which can be any non-negative real number.

The resulting GGM is displayed using a spring layout that places strongly correlated nodes near each other as if the correlation weights are springs pulling nodes together in space. If there are no systematic trends, then the non-zero partial correlations will appear on seemingly random pairs of shots, and the displayed GGM will look like an amorphous *blob* (or *hairball* in graph theory terminology). More visually interesting patterns emerge when there are interesting depth trends, such as a chain for systematic decrease/increase in elements, or clusters for sudden change in chemistry (such as a layer). This automated method identifies compositional depth trends associated with varnish and weathering rinds on laboratory samples [5]; and dust layers and thin sulfate veins on Mars targets [6].

Interactive Gaussian Graphical Models: The GGM gives a quick visual summary of geochemical trends, but to answer specific science questions, we need additional information and control over the machine learning algorithm. We developed an interactive approach to learning a GGM which provides information about which wavelengths or elements that are changing in a depth trend. A scientist can select which wavelengths or elements to include in the analysis which modifies the input data X and updates the learned GGM model [9]. Part of the user interface is shown in Figure 1.

Without the interaction, we can see in the Figure 1 GGM that there is a surface layer on the ChemCam target. In this case, it is a dust layer which we verified by looking at the decrease in abundance of elements associated with martian dust (e.g., Mg [10]) and the increase in abundance with S and Ca. Similar GGM structures are learned from LIBS data on terrestrial lab samples of rocks with rock varnish and weathering

rinds [6]. The specific information about elemental abundance provided by the interactive GGM learning is needed to distinguish between these types of geochemical gradient observations that are produced by very different geological features.

Future of Interactive Machine Learning: Our goal is to help scientists transform data into knowledge by observing how scientists transform data into knowledge. When there are tedious, repetitive steps in the process, then there are opportunities for machine learning to automate the tedium. The interactive machine learning approach opens up new opportunities along two fronts. First, because fully-automated methods are rarely accurate enough for scientific inquiry, interaction allows scientists to fix or fine-tune results that are close but not quite right. Second, by logging user interactions, the machine learning algorithms can improve.

One compelling example of interactive machine learning that we have not yet applied to planetary data, is the segmentation of microscopy images, as shown in Figure 2. The goal is to quantify the size and shape of hundreds of particles in images for quantitative analysis [11]. While it is easy to see the particles, it is tedious for people to draw precise outlines of each particle. Automated image segmentation can find edges precisely, but makes many mistakes in identifying particle boundaries. Our interactive approach starts with the imperfect automated image segmentation, allows the scientist to fix the segmentation mistakes, and logs the entire process. We are developing machine learning algorithms that take the user interactions logs as input data to improve the overall image segmentation.

References: [1] Porter et al. (2013) *Comp. in Sci. & Eng.* [2] Wiens et al. (2012) *Space Sci. Rev.*, 170. [3] Lanza et al. (2015). *Icarus*. [4] Blaney et al. (2014), *JGR*, 119, 2109-2131. [5] Oyen and Lanza. (2015). *LPSC abstract 2940* [6] Oyen and Lanza. (2017). *LPSC abstract 1479*. [7] Koller and Friedman. (2009). *Probabilistic Graphical Models*. [8] Zhao T. et al (2012) *J. Machine Learning Research*. [9] Oyen et al (2016) *Intl. Conf. Artificial Intelligence*. [10] Lasue et al. (2014). *LPSC*, abstract 1224. [11] Porter et al. (2015) *Intl. Symp. Math. Morphology*.

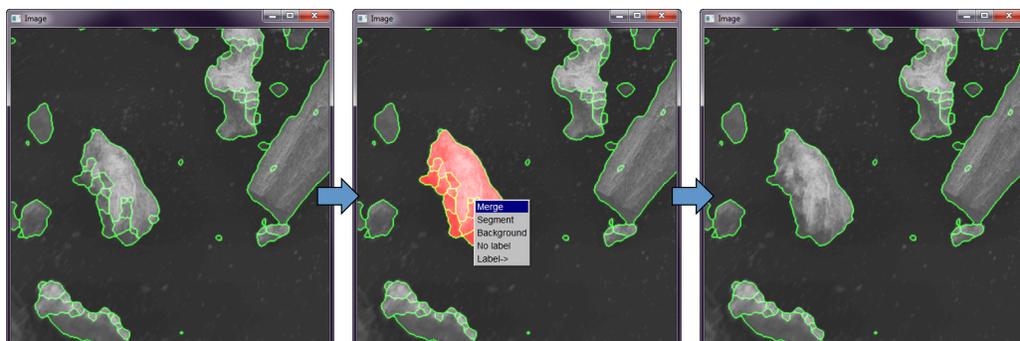


Figure 2 Interactive image segmentation begins with auto segmentation (**left**), user interaction (**middle**), and ends with the corrected segmentation (**right**). Interaction improves the segmentation