# Large Scale Feature Selection and Online Learning
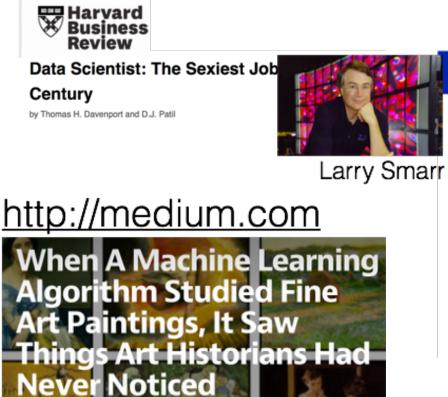
Gregory Ditzler

Dept. of Electrical & Computer Engineering
ditzler@email.arizona.edu
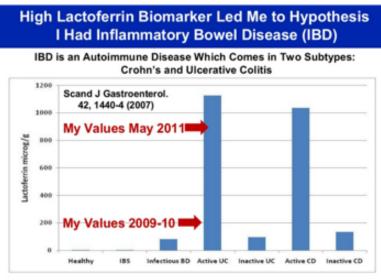
ARIZONA

# Why data science, machine learning and knowledge discovery?

# Challenges & constraints

- There are a lot of data being generated in today's technological climate
  - Pro: some data are useful
  - Con: some data are not useful
  - A little bit like finding the signal in the noise

- Some scenarios request that variables maintain a physical interpretation

- Life sciences and clinical tests

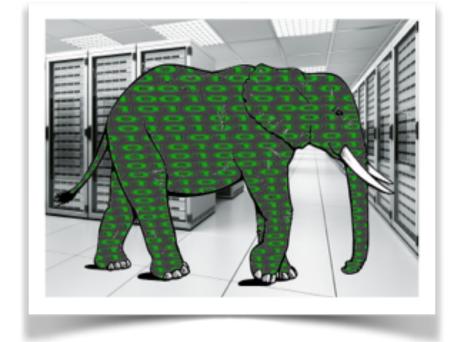# Feature Selection for Large Volumes of Data

# Data aren't small anymore

- Scale of data sets are growing rapidly in the internet era

- What does it mean for data to be "large"?

- The Five V's: volume, velocity, variety, veracity, and value
  - velocity: data arrive in a stream
  - volume: not only number of samples, but the dimensionality
  - value: not of cost, but of importance

- Twenty years of machine learning research has led to wide body of research for detecting value
  - "volume" of twenty years ago is not the volume of today
  - distributed, parallel, and statistically sound

# Feature selection in a nutshell



healthy $y_1$

age
gender
medical tests
$\rightarrow \mathbf{x}_1$

unhealthy $y_2$

age
gender
medical tests
$\rightarrow \mathbf{x}_2 \rightarrow$

Feature Selection

healthy $y_M$

age
gender
medical tests
$\rightarrow \mathbf{x}_M$

$\in \mathbb{R}^K$

Relevance, Weak Relevance
and Irrelevance

$k < K$

$y = \Phi(\mathbf{x})$

Classification

$\mathbf{x}_1'$
$\mathbf{x}_2' \in \mathbb{R}^k$
$\mathbf{x}_M'$

Knowledge
Discovery

Stick figures courtesy of xkcd.com

# Feature selection with an adversary in a nutshell



legitimate $y_1$

malicious $y_2$

legitimate $y_M$

$\mathbf{x}_1$

I know what you're using

$\mathbf{x}_2$

$\mathbf{x}_M$

$f(\mathbf{x})$

legitimate?

$\widehat{y}$

malicious?

Stick figures courtesy of xkcd.com

$\in \mathbb{R}^K$

$$\beta^* = \arg\min_{\beta \in \Phi} \left\{ \underbrace{\ell(\beta, \mathcal{D})}_{\text{model loss}} + \overbrace{\lambda\Omega(\beta)}^{\text{regularization}} + \underbrace{\alpha\Lambda(\beta, \widehat{\mathcal{D}})}_{\text{adversary}} \right\}$$

# Approaches for feature selection

**Wrapper Methods**
- Build a classifier, measure loss, adapt feature set, repeat...
- Easy to over fit
- Too computationally complex

**Embedded Methods**
- Jointly optimize classifier and variable selector parameters
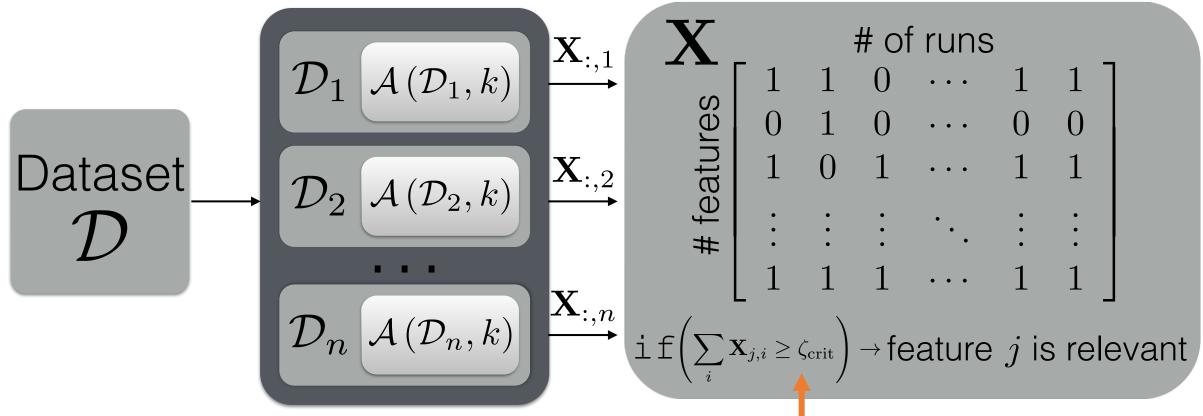- E.g., Linear model with L1 penalization.

**Filter Methods**
- Optimize feature set independent from a classifier
- Fast, but need ways to scale them to large volumes data
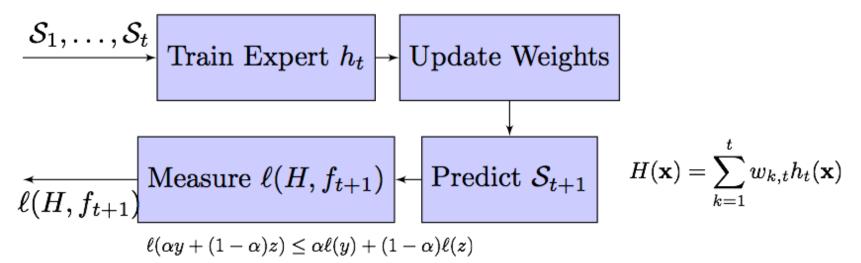
# NPFS in one picture

Map

Reduce & Inference

Dataset $\mathcal{D}$

$\mathcal{D}_1$ $\mathcal{A}(\mathcal{D}_1, k)$

$\mathcal{D}_2$ $\mathcal{A}(\mathcal{D}_2, k)$

$\mathcal{D}_n$ $\mathcal{A}(\mathcal{D}_n, k)$

$\mathbf{X}_{:,1}$

$\mathbf{X}_{:,2}$

$\mathbf{X}_{:,n}$

$\mathbf{X}$    # of runs

$$\text{\# features} \begin{bmatrix} 1 & 1 & 0 & \cdots & 1 & 1 \\ 0 & 1 & 0 & \cdots & 0 & 0 \\ 1 & 0 & 1 & \cdots & 1 & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 1 & 1 & 1 & \cdots & 1 & 1 \end{bmatrix}$$

$\mathrm{if}\left( \sum_i \mathbf{X}_{j,i} \geq \zeta_{\mathrm{crit}} \right) \rightarrow$ feature $j$ is relevant

Neyman-Pearson hypothesis test for relevance

# Streaming Data and Learning Online

# Learning in Nonstationary Environments

$$\mathcal{S}_1, \ldots, \mathcal{S}_t$$

Train Expert $h_t$ → Update Weights

Measure $\ell(H, f_{t+1})$ ← Predict $\mathcal{S}_{t+1}$

$\ell(H, f_{t+1})$

$$H(\mathbf{x}) = \sum_{k=1}^{t} w_{k,t} h_t(\mathbf{x})$$

$$\ell(\alpha y + (1-\alpha)z) \le \alpha \ell(y) + (1-\alpha)\ell(z)$$

- The Five V's: volume, velocity, variety, veracity, and value
- Traditional Learning Paradigm: training and testing data are sampled from the same probability distribution
  - what if that is not the situation?
- What should we expect the loss to look like when predicting on data from an unknown distribution?

# Discussion Points

"There are two types of machine learning practitioners: (1) those that generalized from limited data"

# Challenges Moving Forward

- Advanced Frameworks for Big Data Subset Selection
  - IEEE CIM recently published a special issue on Big Data and the curse of big dimensionality
  - Millions of features & beyond while being mindful of veracity
  - Interpretability, Visualization and Real-Time
  - Migrating from a batch-based learning setting to a pure online setting

- Calibrated Prediction for Domain Adaptation in Time-Series
  - Tuning model parameters in changing domains using unlabeled data, and accessing the stability of predictions in uncertain environments

- Learning New Classes without Re-training

Thanks for Listening!

Acknowledgements