January 28, 2002 , COMS 6998-01   Advanced Machine Learning

Prof. Tony Jebara , CEPSR 605    jebara @ cs.columbia.edu
http \\ www.cs.columbia.edu \ ~jebara \ 6998-01
Office Hours: Tuesday 2-4, Thursday 11am-12am, Or Appt

CS Dept
Area: ML
for modeling
people, behavior,
medical, visual

Text book Online, Duda Hart & Stork , Bishop
          Email me for permission

Handout Outline (web page) , Grading Homework & Project
Pre-reqs: Linear Algebra, Calculus
          Intro. Machine Learning  or Stats

Class List :  names , emails,  listener  or  for credit

Course Structure:    Seminar-like,   hands-on,
          interactive, MATLAB,  get  your hands dirty,
          researchy topics,  new and controversial  ideas
          as well as formal.   Email-list: discussion-group.
          Learn  a  set  of  tools  for  many applications
                & research.

ML:   fields are separating,   hard to keep up  Jack of all trades;
      Physics, Math, Statistics, Economics, OR, Neuroscience, Psych, Biology

ML: model complex and non-deterministic systems

Physics $E=MC^2$, many other domains: unknown eqns. on vars

- partial observations
- unknown / uncertain / incomplete models
- high dimensions   100 000 vars, not 3
- noise
- complex   (non linear)
- stochastic
- refine models with data / real measurements and observations

Applications:
- Speech Rec  ⛓ HMMs
- Computer Vision   (face rec, digits, ▦ MRFs Super-res.)

all share
common tools
and theory

- Time Series Prediction (weather, financial...)
- Genomics  (DNA, micro-arrays, SVMs, splice sites)
- NLP & Parsing   HMMs, CRFs  ⛓, Whizbang
- Text classif. & IR ( documents, spam, TSVMs)
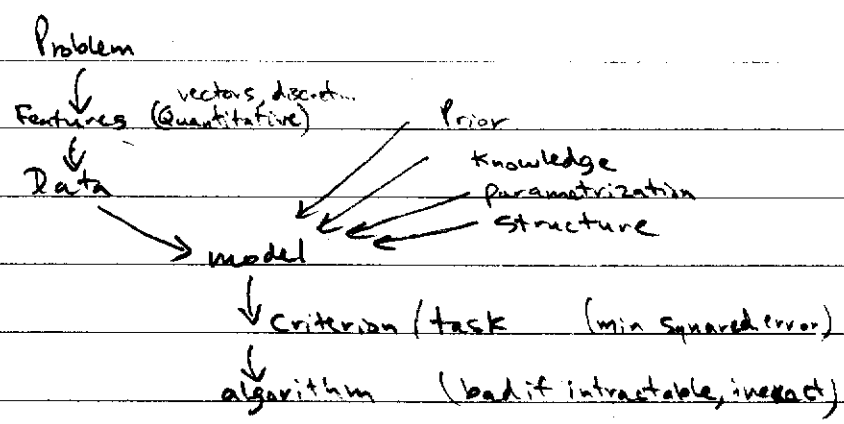- Medical: QMR-DT    600 diseases / 4000 symptoms

History:
1917: Karel Capek - Robot
1943: McCulloch & Pitts (Bio, Neuron)                    Statistics
1947: Norbert Wiener, Cybernetics, multi-fields
1949: Shannon
1950's: Minsky, Allen Newell, Herbert Simon, John McCarthy          ~70 HMM Baum
         Symbolic AI, Logic, INCONSISTENCY, Rule-based           1977 EM (Demp)
1957: Rosenblatt Perceptron
                     kill AI, 12-15yrs
1969: Minsky & Papert: Perceptron linear ✗        Graphical Models 1980
                                                   Lauritzen, Pearl
1974: Werbos, PhD Backprop, Nonlinear
1986: Rumelhart & McLennan, MLP, Conjugate

1980's: NN, RNN, Genetic Alg., Fuzzy Logic, Black Boxes.

1990's: Bayesian & Statistical & Structure & Priors

Graphical Models: EM, KF, HMM, Sig. Belief Nets, MRFs

Course → SVMs, Comput. Learning Theory, Boosting

Course Outline:  Slow & faster, email me on pace
old topics + new research : EM - 3 views

General ML Recipe:        Problem
                            ↓
                Features (Quantitative)   vectors, discrt...      Prior
                            ↓                                   Knowledge
                Data                                            parametrization
                            ↘                                   Structure
                            → model
                                ↓
                                Criterion (task    (min squared error)
                                ↓
                                algorithm    (bad if intractable, inexact)

Tasks:  – modeling & density estimation   (compression, coding, analysis)
        – clustering
Overlap – classification < binary / multiclass    $x_1, ... x_T$   $y_1, ..., y_T$
        – regression
        – structure learning
        – feature selection, subspace
        – transduction
        – anomaly detection

Discriminative vs. Generative:    $p(x, y, ....)$       analytic
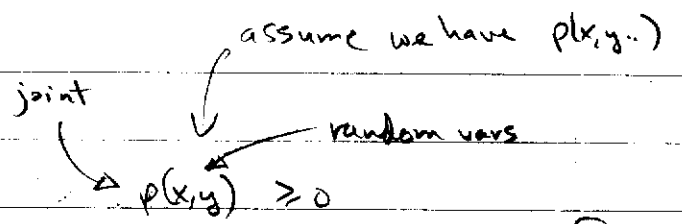                                                        pdf over
                                                        all variables
                                                        in system

        general yet suboptimal ←              →  get whatever
                                                 by manipulating this.

Probability Theory Review:    $\xrightarrow{\text{joint}}$ $p(x,y) \geq 0$    $\xleftarrow{\text{random vars}}$

(useful properties)    $\sum_{x,y} p(x,y) = 1 \quad \longrightarrow \quad \sum_{x,y} p(x,y) = 1$

assume we have $p(x,y..)$

marginalize:    $\sum_{y} p(x,y) = p(x)$

condition:    $p(x|y) = \dfrac{p(x,y)}{p(y)}$

bayes rule:    $p(x,y) = \dfrac{p(y|x) p(x)}{p(y)}$

$x \perp\!\!\!\perp y$ : independent $\Longrightarrow$ $p(x,y) = p(x)p(y)$  OR  $p(x|y) = p(x)$

$x \perp\!\!\!\perp z \,|\, y$ : cond indep $\Longrightarrow$ $p(x|z,y) = p(x|y)$

$\xrightarrow{\text{BUT}}$ $p(x|z) \neq p(x)$

classifier: $p(y|x)$    $\hat{y} = \underset{y}{\arg\max} \{ p(y=0|\hat{x}) , p(y=1|\hat{x}) ... \}$

regression: $p(y|x)$    $\hat{y} = \underset{q}{\arg\max} \; p(y=q|\hat{x})$

$\hat{y} = \int y \, p(y|\hat{x}) \, dy$

anomaly detection:    $p(x) \geq$ threshold $\Longrightarrow$ true, else anomaly

Graphical Model Rep:    nodes $\rightarrow$ variables

indep. $\Rightarrow$ no link    ⓧ  ⓨ

depend $\equiv$ link    ⓧ —→ ⓨ
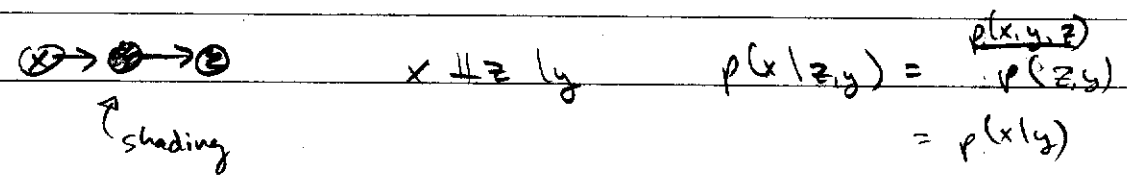
acyclical    arrow $\equiv$ parent-child ($\sim$ causal)

$p(x_1,...,x_n) = \prod_i p(x_i | pa_i) = \prod_i p(x_i | \pi_i)$
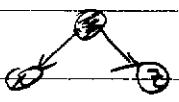
↑ efficiency, tables, later on...

$k^n$ vars    $\times$ ▦

chain of events

ⓧ→ⓨ→ⓩ    $p(x) \, p(y|x) \, p(z|y)$

$\begin{bmatrix} X = \text{trip} \\ Y = \text{fall down stairs} \\ Z = \text{bruise} \end{bmatrix}$

ⓧ→⬤→ⓩ    $x \perp\!\!\!\perp z \,|\, y$    $p(x|z,y) = \dfrac{p(x,y,z)}{p(z,y)}$

↑ shading    $= p(x|y)$

One cause* with 2 effects

$$Y \rightarrow X, \quad Y \rightarrow Z$$

$$\begin{cases} y = \text{flu} \\ x = \text{sore throat} \\ z = \text{temperature} \end{cases}$$

$$x \perp\!\!\!\perp z \mid y$$

2 causes*, 1 effect
explaining away

$$x \perp\!\!\!\perp z$$

$$x \not\perp\!\!\!\perp z \mid y$$

$$\begin{cases} x = \text{rain} \\ y = \text{wet driveway} \\ z = \text{car oil leak} \end{cases}$$

— How to get $P(x, y, ...)$? Parametric model $P_\theta(x, y, ... | \theta)$
Estimation Two Schools: Bayesian & Frequentist
Both have own advantages, we focus on Bayesian

Frequentist: Classical, objective, no ~~priors~~
Some things are not distributions
Can talk about p(event) if never get data.
$p(\theta)$ , only 1 true model, not R.V.
$p_\theta(x, y)$     $p(x, y | \theta)$
Can't say prob. a coin will be 50% heads w/o observ.
plug in Single $\theta$ is generated with estimator, MLE
min variance, unbiased estimator

Bayesian: Subjective, pdf on anything (i.e. uncertainty)
statisticians
even on deterministic values, like speed of light
pure

Bayes rule:    $$p(\theta | X) = \frac{p(X | \theta) \, p(\theta)}{p(X)}$$

You pick it
simplicity
Ockham
1280-1349
mathematical
convenience

distrib over $\hat\theta$     posterior $= \dfrac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$

Bayes 1702-1761