# Midterm Review

Machine Learning 10-701

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

March 1, 2010

See practice exams at:
http://www.cs.cmu.edu/~tom/10601_sp09/601-sp09-midterm-solutions.pdf
http://select.cs.cmu.edu/class/10701-F09/exams.html

**Midterm is open book, open notes, NO computers**

**Covers all material presented up through today's class.**

# Some Topics We've Covered

## Decision trees

entropy, overfitting

## Probability basics

rv's, manipulating probabilities,
Bayes rule, MLE, MAP,
conditional indep.

## Instance-based learning

nearest nbr., density estimation,
Bayes optimal classifier

## Naïve Bayes

conditional indep, # of parameters
to estimate,

## Logistic regression

form of $P(Y|X)$ implied by N. Bayes,
generative vs. discriminative

## Linear Regression

minimizing sum sq. error ~ MLE
regularization ~ MAP, non-linear

## Neural Networks

gradient descent,
learning hidden representations

## Model Selection

overfitting, bias-variance

## Clustering

k-means, mixture Gaussians, EM

## Hidden Markov Models

time series model, backward-forward

## Bayesian Networks

factored representation of joint
distribution, encoding conditional
independence assumptions

| | representation of P(Y\|X) | decision surface | optimization objective | convergence guarantee? | other assumptions? |
|---|---|---|---|---|---|
| Naïve Bayes | | | | | |
| Logistic Regr. | | | | | |
| Linear Regr. | | | | | |
| Neural net | | | | | |
| Dec. Tree | | | | | |
| Gaussian Mixture model | | | | | |
| HMM | | | | | |
| Bayes Net | | | | | |
| kNN | | | | | |

# Four Fundamentals for ML

1. Learning is an optimization problem

2. Learning is a parameter estimation problem

3. Error arises from three sources

4. Practical learning requires modeling assumptions, such as …
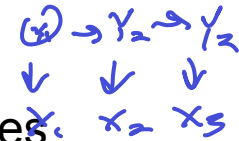
# Learning is an optimization problem

- – many algorithms are best understood as optimization algs
- – what objective do they optimize, and how?
- – naïve Bayes?  logistic regression?  linear regression?

suppose want $\boxed{f(x) = y}$, $P(y|x)$

# Learning is parameter estimation  HMM

$\theta$ defines learned $f$.

$$\text{(s)} \to Y_2 \to Y_3$$
$$\downarrow \quad \downarrow \quad \downarrow$$
$$x_1 \quad x_2 \quad x_3$$

- the more training data, the more accurate the estimates
- to measure accuracy of learned model, we must use test (not $P(data|\theta)$ train) data

$$= \prod_k P(x_1 \ldots x_3^k | \theta)$$

- cross validation

$N(0, \sigma)$

$$N\,Bayes \quad P(y|x_1 \ldots x_n) = \boxed{\prod_i P(x_i | y)} \cdot \boxed{P(y)} \quad \text{if assume}$$
$$Y = f(x, \theta) + \epsilon$$

$$\hat{P}(y|x, \theta)$$
$$\underset{K\,params}{}$$

Lin $\quad \underset{\theta}{argmax} \prod_k P(y^k | x^k, \theta)$

$$\text{trainy examps:} \quad x^1 y^1, \; x^2 y^2 \ldots x^k y^k \quad \text{regr.} \overset{\theta}{\to} \hat{\theta} = argmin \sum_k (y^k - \hat{f}(x^k, \theta))^2$$

$$\underset{=}{P(data|\theta)} \; data\;likelihood\,(\theta) = \prod_{k=1}^{K} P(x^k y^k | \theta)$$

MCLE

Cond. likelihood $= \prod_k P(y^k | x^k, \theta)$

$$MLE =$$
$$\underset{\theta}{argmax}\; data\,likelihood(\theta) = \prod^{\boxed{K}} P(x^k | y^k, \theta) \; P(y^k | \theta)$$

$$MAP \neq \underset{\theta}{argmax} \boxed{P(\theta | data) = \frac{P(data|\theta)\;P(\theta)}{P(data)}}$$
$$k=1$$

# Error arises from three sources

– Bayes optimal error, bias$^2$, variance

learning $f_\theta(x)$ means picking $\theta$

r.v. because sample of training data is drawn randomly from $P(x,y)$
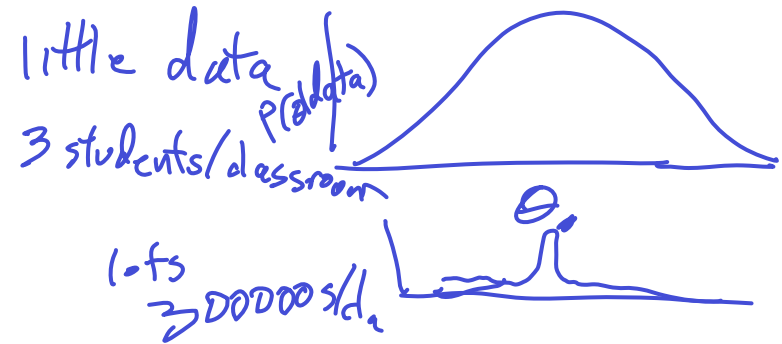
Inescapable error even if know perfectly $P(Y|X)$

$E[\hat\theta] \neq \theta$

$E[\hat\theta] - \theta$

$E[\hat\theta - E[\hat\theta)]$

Usually simpler models exhibit higher bias, lower variance

more data $\Rightarrow$ lower variance

little data
3 students/classroom

$P(\theta|data)$

lots
300000 s/d

$\theta$

# Bias and Variance

$= \hat{\theta}$ estimate

given some estimator Y for some parameter θ, we note
  Y is a random variable (why?)

the <u>bias</u> of estimator Y : $E[Y] - \theta$ — if Y is <u>un</u>biased then $E[Y] = \theta$

the <u>variance</u> of estimator Y : $E[(Y - E[Y])^2]$

expectation is over different
draws of training data

consider when

- θ is the probability of "heads" for my coin
- Y = proportion of heads observed from 3 flips

# Practical learning requires making assumptions

- Why?
- form of the f:X → Y, or P(Y|X), or P(…) to be learned
- priors on parameters → MAP, regularization
- Conditional independence → Naive Bayes, Bayes nets

# Four Fundamentals for ML

1. ## Learning is an optimization problem
   - many algorithms are best understood as optimization algs
   - what objective do they optimize, and how?

2. ## Learning is a parameter estimation problem
   - the more training data, the more accurate the estimates
   - MLE, MAP, M(Conditional)LE, …
   - to measure accuracy of learned model, we must use test (not train) data

3. ## Error arises from three sources
   - Bayes optimal error, bias, variance

4. ## Practical learning requires modeling assumptions
   - Why?
   - form of the f:X → Y, or P(Y|X) to be learned
   - priors on parameters: MAP, regularization
   - Conditional independence: Naive Bayes, Bayes nets, HMM's