# CSE 573: Artificial Intelligence
## Spring 2012

Structure Learning, EM, Cotraining

Dan Weld

Slides adapted from Carlos Guestrin, Krzysztof Gajos, Dan Klein, Stuart Russell, Andrew Moore & Luke Zettlemoyer

---

## Some Typical Biases

▪ Occam's razor

We are to admit no more causes of natural things than such as are both true and sufficient to explain their appearances
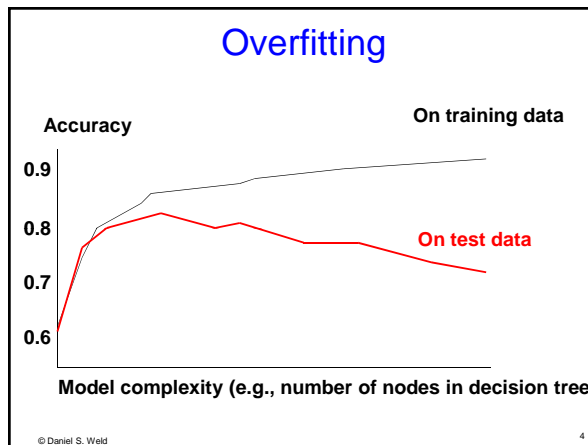– William of Ockham (1288-1348)



© Daniel S. Weld                                                   2
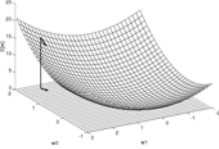
---

## Some Typical Biases

▪ Occam's razor
▪ MDL – Minimum description length
▪ Concepts can be approximated by
  ... **conjunctions** of predicates,
  ... **linear** functions
  ... **short** decision trees
▪ Maximal conditional independence
▪ Minimum cross-validation error
▪ Minimum number of features
▪ Etc..

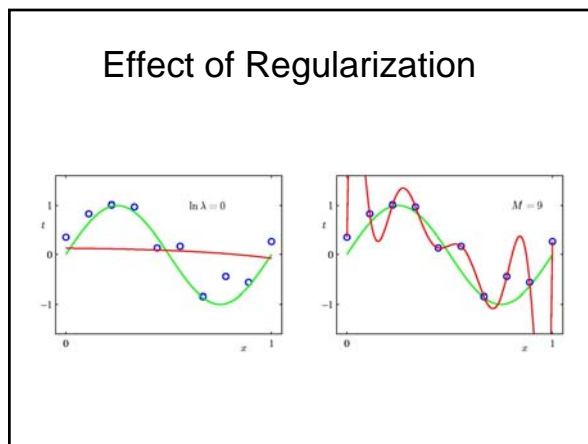© Daniel S. Weld                                                   3

---

## Overfitting



Accuracy

On training data

On test data

0.9

0.8

0.7

0.6

Model complexity (e.g., number of nodes in decision tree)

© Daniel S. Weld                                                   4

---

## Learning as Optimization

▪ Methods
  ▪ Closed form
  ▪ Greedy search
  ▪ Gradient ascent
▪ Loss Function (preference bias)
  ▪ Minimize *loss* over training data (test data)
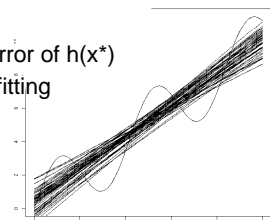  ▪ Loss(h,data) = error(h, data) + complexity(h)

Regularization term    E.g., $\lambda\,||w||^2$
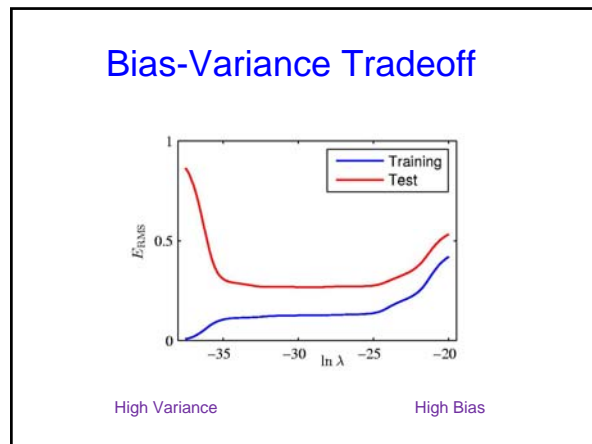
---

## Effect of Regularization



---

## Bias / Variance Tradeoff

- Variance: **E[ (h(x*) – h(x*))² ]**
  How much h(x*) varies between training sets
  Reducing variance risks underfitting

- Bias: **[h(x*) – f(x*)]**
  Describes the *average* error of h(x*)
  Reducing bias risks overfitting

Note: **inductive bias** *vs* **estimator bias**

Slide from T Dietterich

## Bias-Variance Tradeoff
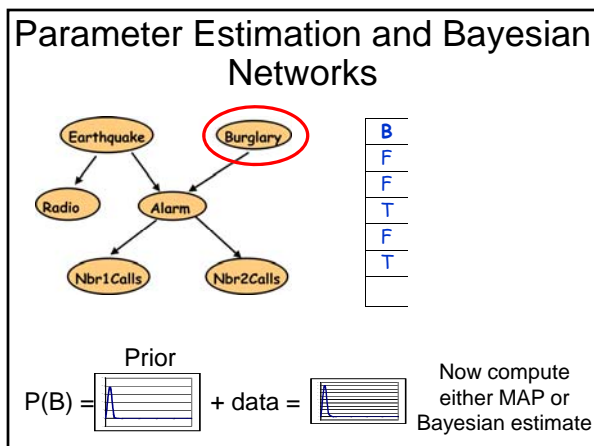
High Variance                     High Bias

## Topics

- Learning Parameters for a Bayesian Network
  - Fully observable
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks
- Cool Stuff
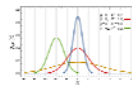  - Learning Ensembles
  - Cotraining

© Daniel S. Weld

## Summary

| | Prior | Hypothesis |
|---|---|---|
| Maximum Likelihood Estimate | Uniform | The most likely |
| Maximum A Posteriori Estimate | Any | The most likely |
| Bayesian Estimate | Any | Weighted combination |

Minimizes error
Great when data is scarce
Potentially much harder to compute

Beta & Dirichlet

## Parameter Estimation and Bayesian Networks

| B |
|---|
| F |
| F |
| T |
| F |
| T |

Prior

P(B) = [plot] + data = [plot]   Now compute either MAP or Bayesian estimate

## Learning with Continuous Variables

Earthquake

Pr(E=x)
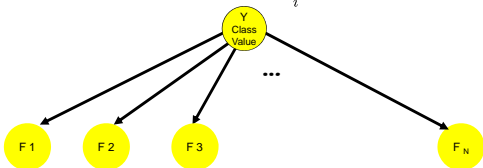mean: μ = **?**
variance: σ = **?**

$$\widehat{\mu}_{MLE} \;=\; \frac{1}{N}\sum_{i=1}^{N} x_i$$

$$\widehat{\sigma}^2_{MLE} \;=\; \frac{1}{N}\sum_{i=1}^{N} (x_i - \widehat{\mu})^2$$

© Daniel S. Weld

## A Popular Structure: Naïve Bayes

$$P(\mathsf{Y}, \mathsf{F}_1 \ldots \mathsf{F}_n) = P(\mathsf{Y}) \prod_i P(\mathsf{F}_i | \mathsf{Y})$$
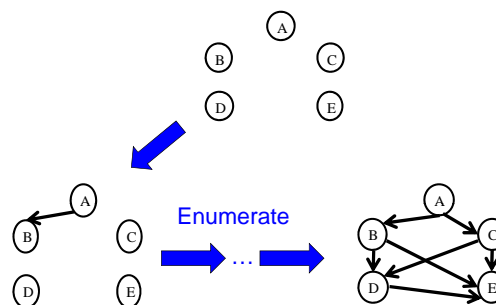


Assume that features are conditionally independent given class variable
Works surprisingly well for **classification** (predicting the right class)
But forces probabilities towards 0 and 1

---

What if we ***don't*** know structure?

---

## Learning The Structure of Bayesian Networks

- Search thru the space…
  - of possible network structures!
  - (for now still assume can observe all values)
- For each structure, learn parameters
  - As just shown…
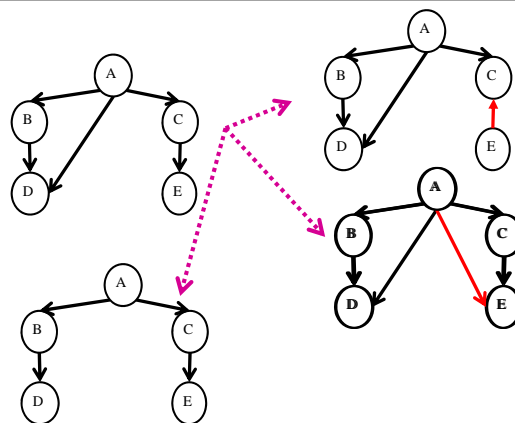- Pick the one that fits observed data best
  - Calculate P(data)

---



**Two problems:**
- Fully connected graph will be most probable
- Exponential number of structures

---

## Learning The Structure of Bayesian Networks

- Search thru the space…
  - of possible network structures!
- For each structure, learn parameters
  - As just shown…
- Pick the one that fits observed data best
  - Calculate P(data)

**Two problems:**
- Fully connected will be most probable
  - Add penalty term (regularization) ∝ model complexity
- Exponential number of structures
  - Local search

---

## Score Functions

- Bayesian Information Criterion (BIC)
  - P(D | BN) – penalty
  - Penalty = ½ (# parameters) Log (# data points)

- MAP score
  - P(BN | D) = P(D | BN) P(BN)
  - P(BN) must decay exponentially with # of parameters for this to work well

- Loss(h,data) = error(h, data) + complexity(h)

© Daniel S. Weld
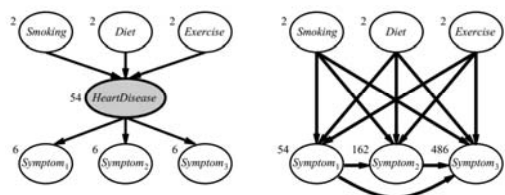
20

## Topics

- Learning Parameters for a Bayesian Network
  - Fully observable
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks
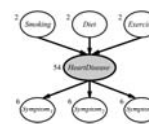- Cool Stuff
  - Learning Ensembles
  - Cotraining

© Daniel S. Weld

## Why Learn Hidden Variables?



## Chicken & Egg Problem

- If we knew whether patient had disease
  - It would be easy to learn CPTs
  - But we can't observe states, so we don't!



- If we knew CPTs
  - It would be easy to predict if patient had disease
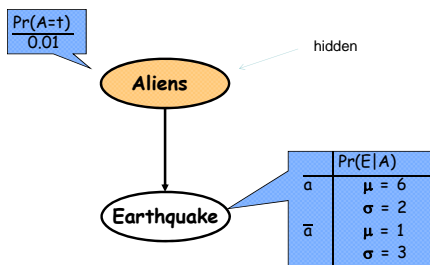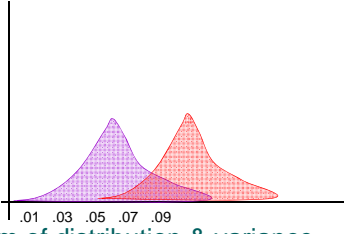  - But we don't, so we can't!

Slide by Daniel S. Weld

23



## Continuous Variables



Pr(A=t)
0.01

Aliens

hidden

Earthquake

| | Pr(E|A) |
|---|---|
| a | μ = 6 |
| | σ = 2 |
| ā | μ = 1 |
| | σ = 3 |

© Daniel S. Weld

4

## Simplest Version

- Mixture of two distributions



- Know: form of distribution & variance, $\sigma = 1$
- Just need *mean* of each distribution

.01 .03 .05 .07 .09

26

## Input Looks Like



.01 .03 .05 .07 .09

27

## We Want to Predict



**?**

.01 .03 .05 .07 .09

28

## Chicken & Egg

Note that coloring instances would be easy *if* we knew Gaussian parameters….



.01 .03 .05 .07 .09

29

## Chicken & Egg

And finding the Gaussians would be easy *if* we knew the coloring



.01 .03 .05 .07 .09

30

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly: set $\theta_1 = ?$; $\theta_2 = ?$



.01 .03 .05 .07 .09

31

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable

.01  .03  .05  .07  .09

32

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
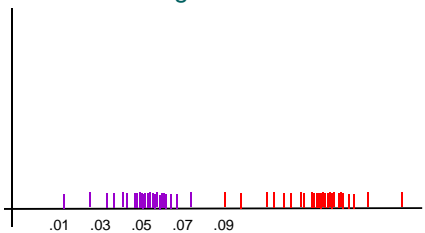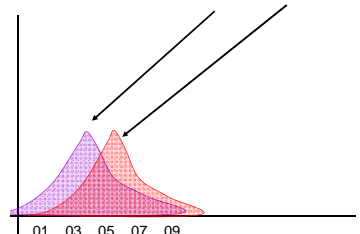  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable

.01  .03  .05  .07  .09

33

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable

**[M step]** Treating each instance as *fractionally* having **both** values compute the new parameter values

.01  .03  .05  .07  .09

34

## ML Mean of Single Gaussian

$$U_{ml} = \text{argmin}_u \sum_i (x_i - u)^2$$

.01  .03  .05  .07  .09

35

## Expectation Maximization (EM)

- 

**[M step]** Treating each instance as fractionally having **both** values compute the new parameter values

.01  .03  .05  .07  .09

36

## Expectation Maximization (EM)

- [E step] Compute probability of instance having each possible value of the hidden variable

.01  .03  .05  .07  .09

37

## Expectation Maximization (EM)

- [E step] Compute probability of instance having each possible value of the hidden variable

  [M step] Treating each instance as fractionally having both values compute the new parameter values

.01  .03  .05  .07  .09

38

## Expectation Maximization (EM)

- [E step] Compute probability of instance having each possible value of the hidden variable

  [M step] Treating each instance as fractionally having both values compute the new parameter values

.01  .03  .05  .07  .09

39

## EM

- Works for multiple hidden variables
  & other parametric forms
  - E.g., Baum-Welch algorithm for HMMs

- Optimality?
- Complexity?

- Search?

40

## Topics

- Learning Parameters for a Bayesian Network
  - Fully observable
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks
- Cool Stuff
  - Learning Ensembles
  - Cotraining

## Ensembles of Classifiers

- Traditional approach: Use one classifier
- Can one do better?
- Approaches:
  - Cross-validated committees
  - Bagging
  - Boosting
  - Stacking

## Ensembles of Classifiers

- Assume
  - Errors are independent (suppose 30% error)
  - Majority vote
- Probability that majority is wrong…

  = area under binomial distribution

Prob 0.2

0.1

Ensemble of 21 classifiers

Number of classifiers in error

- If individual area is 0.3
- Area under curve for ≥11 wrong is 0.026
- Order of magnitude improvement!

## Constructing Ensembles
## Cross-validated committees

- Partition examples into *k* disjoint equiv classes
- Now create *k* training sets
  - Each set is union of all equiv classes *except one*
  - So each set has (k-1)/k of the original training data

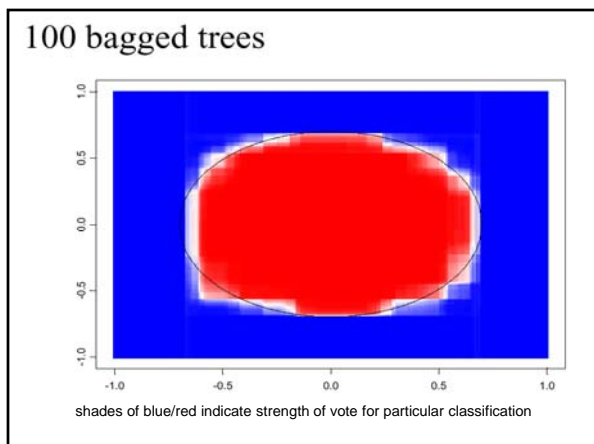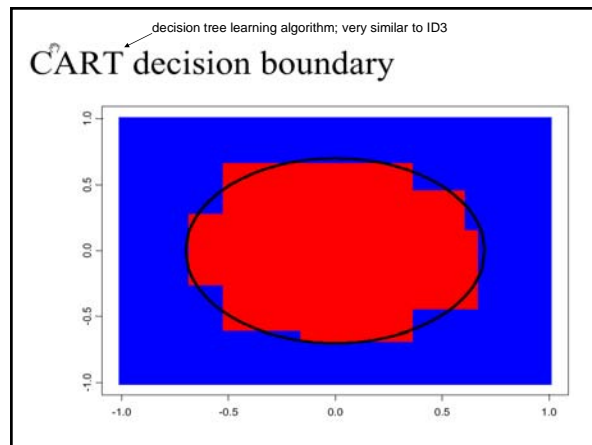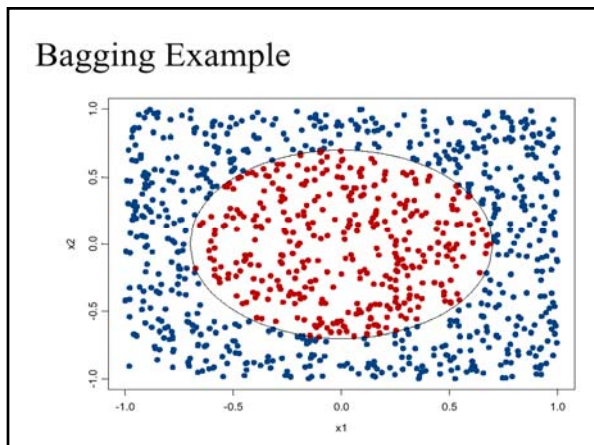- Now train a classifier on each set



44                                        © Daniel S. Weld

## Ensemble Construction II
## Bagging

- Generate k sets of training examples
- For each set
  - Draw m examples randomly (with replacement)
  - From the original set of m examples
- Each training set corresponds to
  - 63.2% of original (+ duplicates)
- Now train classifier on each set
- Intuition: Sampling helps algorithm become more robust to noise/outliers in the data

45                                        © Daniel S. Weld

## Bagging Example



## CART decision boundary

decision tree learning algorithm; very similar to ID3



## 100 bagged trees



shades of blue/red indicate strength of vote for particular classification

## Boosting    [Schapire, 1989]

- Idea: run weak learner multiple times on (reweighted!) training data; weight learned classifiers $\propto$ their accuracy

- On each iteration *t*:
  - Learn a hypothesis, $h_t$, using distribution to weight examples
  - Compute a strength for this hypothesis – $\alpha_t$
  - Reweight training examples by how well they were classified

- Final classifier:

$$h(x) = \text{sign}\left(\sum_i \alpha_i h_i(x)\right)$$

- **Practically useful**
- **Theoretically interesting**

## Bagging vs Boosting

Bagged Decision Rule

Boosted Decision Rule

52

## Ensemble Creation IV
## Stacking

- Train several base learners
- Next train meta-learner
  - Learns when base learners are right / wrong
  - Now meta learner arbitrates

Example → Decision Tree → Naive Bayes → Neural Net → Meta−Learner → Prediction

Train using cross validated committees
- Meta-L inputs = base learner predictions
- Training examples = 'test set' from cross validation

53

© Daniel S. Weld

## Topics

- Learning Parameters for a Bayesian Network
  - Fully observable
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks
- Cool Stuff
  - Learning Ensembles
  - Semi-supervised learning (Cotraining)

© Daniel S. Weld