

Hidden Values
Expectation Maximization
Mixture Models

Admin

- Marked Assn 4 due now
- Assn 5 due Wednesday
- Project Proposal due coming Monday
- AG out tonight, due next Wednesday
- Midterm 2 weeks from today (study guide coming soon)

Question From Last time:

Why does bootstrap select ~63% for large N ?

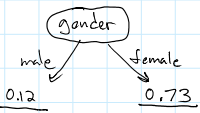
$$\begin{aligned}
 p(i \text{ selected at least once}) &= 1 - p(\text{not selected } N \text{ times}) \\
 &= 1 - (1 - \frac{1}{n})^n \\
 &= 1 - \frac{1}{e} \\
 &\approx 63\%
 \end{aligned}$$

Example of Decision Tree

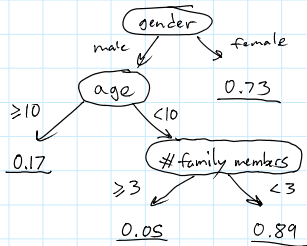
(same as on wikipedia)

	gender	Age	# Family Members	
$x =$	Male	33	5	$y =$ 'died'
	Female	10	1	'lived'
	\vdots	\vdots	\vdots	\vdots

Decision Stump *looks at one variable*

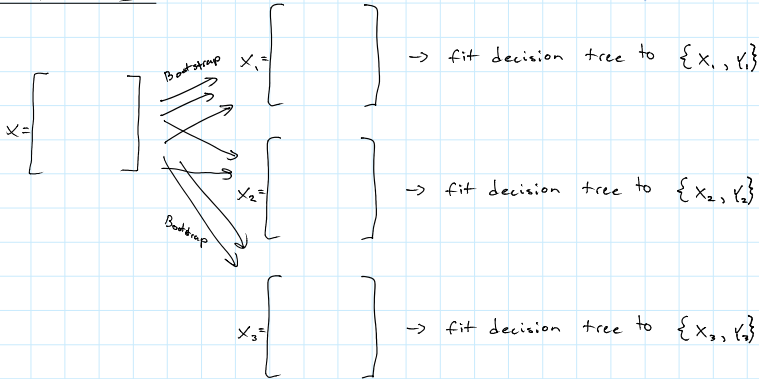


Decision Tree



Random Forests

make a bunch of data sets independently



Each decision node is "Random Decision Stump"

- picks random features (subset of features)
- use info gain
- you can prune with "left-out" samples

$$x = \begin{bmatrix} x_1 \\ \text{[shaded area]} \end{bmatrix} \left. \vphantom{x} \right\} \sim 63\% \text{ of original data}$$

Hidden Values

- learning when some values are unobserved, missing, hidden, latent

	Gender	Age	# Family Members	
$x =$	Male	33	5	$y =$ 'died'
	Female	10	1	'lived'
	Female	?	2	'died'
	Male	22	0	?

Semi Supervised Learning

Idea: getting labels is expensive, getting unlabeled data is cheap

- can we train on $\{X_L, Y_L\}$ and $\{X_U\}$

$$X_L = \begin{bmatrix} \\ \\ \\ \end{bmatrix} \quad Y_L = \begin{bmatrix} \\ \\ \\ \end{bmatrix}$$

$$X_U = \begin{bmatrix} \\ \\ \vdots \\ \end{bmatrix} \quad Y_U = ?$$

(Very large)

Information inequality:
more data can't hurt!

"Missing at Random" (MAR)

- The fact that it is missing does not depend on the missing value.

- E.g. Digit Classification

$$\boxed{2} \rightarrow 2.$$

- missing random pixels: MAR

$$\boxed{3} \rightarrow 3.$$

- hide the labels of all the "2"s examples (not MAR)

- hide the top half of every digit



- If not MAR, you need to model WHY data is missing.

Approach #1

1. Imputation: replace ?'s with most likely value

ie. guess the age if its missing

2. Fit model with "imputed values"

3. Can impute missing data with new, fitted model to try and be better

"hard-EM"

Probabilistic Approach

Notation X : observed variables H : hidden variables

$$P(X) = \sum_h p(X, H=h) \quad (\text{integral if } h \text{ is a continuous random variable})$$

Eg. SSL

$$P(\bar{y}_L, X_L, X_U) = \prod_{i=1}^N p(y_i, x_i) \prod_{j=1}^T \left(\sum_{x_j} p(y_j, x_j) \right)$$

Problem: We assume $-\log P(X, H)$ is "nice" (closed-form, convex)

Maximize $p(x)$:

$$\log(p(x)) = \log \left(\sum_h p(X, H=h) \right)$$

if probability is exponential,
we have familiar log-sum-exp

Minimizing: $-\log(p(x)) = -\log \left(\sum_h p(X, H=h) \right)$

not convex

-problem is the summation inside log

side note:

$$\log(1 + \exp(w^T x))$$

$$\log(\exp(l) + \exp(w^T x)) \quad \text{"convex"}$$

Expectation Maximization

Local Optimizer when $\log p(X, H)$ is "nice" with parameters θ

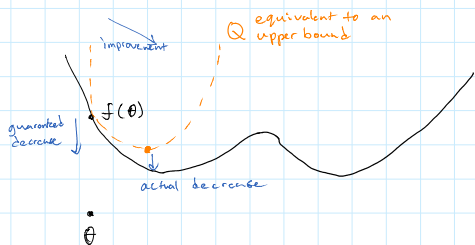
Problem: $\max_{\theta} p(X|\theta)$

Iterations θ^t

$$E_x[f(x)] = \sum_x p(x) f(x) \quad E_{x|y}[f(x)] = \sum_x p(x|y) f(x)$$

"E-step": Define $Q(\theta|\theta^t) = E_{H|X, \theta^t} [\log p(X, H|\theta)]$

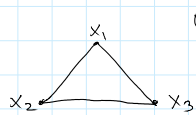
"M-step": $\theta^{t+1} = \arg\max_{\theta} Q(\theta|\theta^t) = \sum_h \underbrace{p(h|X, \theta^t)}_{(\text{fixed}) \alpha_h^t} \underbrace{\log p(X, h|\theta)}_{\text{'nice'}}$



Theorem

$$\log(p(x|\theta^{t+1})) - \log p(x|\theta^t) \geq Q(\theta^{t+1}|\theta^t) - Q(\theta^t|\theta^t)$$

"Convex Combination"



$$\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n$$

$$\alpha_i \geq 0 \quad \sum_i \alpha_i = 1$$

If f is convex,

$$f\left(\sum_i \alpha_i x_i\right) \leq \sum_i \alpha_i f(x_i)$$

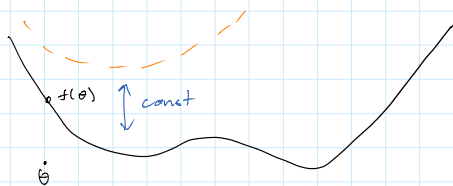
$$\begin{aligned} -\log p(x|\theta) &= -\log \sum_h p(x, h|\theta) \\ &= -\log \left(\sum_h \alpha_h \frac{p(x, h|\theta)}{\alpha_h} \right) \end{aligned}$$

multiply & divide by α_h



$$\leq -\sum_h \alpha_h \log \frac{p(x, h|\theta)}{\alpha_h}$$

$$\begin{aligned} &= -\sum_h \alpha_h \log p(x, h|\theta) + \underbrace{\sum_h \alpha_h \log \alpha_h}_{-H(\alpha) \text{ negative entropy}} \\ &= -Q(\theta|\theta^*) + \text{const} \end{aligned}$$



$$\frac{p(x|H, \theta)}{p(H|x, \theta)} = \frac{p(x, H|\theta)}{p(H|x, \theta)}$$

$$-\log(p(x|\theta)) = -\log(p(x, H|\theta)) + \log(p(H|x, \theta))$$

take expectation:

$$\begin{aligned} \sum_h \alpha_h (-\log p(x|\theta^*)) &= \sum_h \alpha_h (-\log p(x, H|\theta^*) + \log p(H|x, \theta^*)) \\ &= \sum_h \left[\alpha_h \log p(x, H|\theta^*) + \underbrace{\alpha_h \log p(H|x, \theta^*)}_{\alpha_h} \right] \\ &= Q(\theta^*|\theta^*) + \text{const} \end{aligned}$$

"it does touch the point"



Mixture Models

- Recall fitting Gaussians:



- Want to fit probabilistic model

- don't have labels
- but it would be easy if we did!
- motivates mixture models

Let $Z_i \in \{1, 2, \dots, K\}$

$$p_k(x_i | \theta) = p(x_i | z_i = k, \theta)$$

"Gaussian"

Z_i is a "latent" variable never observe, but would be nice to know

$$p(x_i | \theta) = \sum_{k=1}^K p(x_i, z_i = k | \theta)$$

\sum_k for full data set

$$= \sum_{k=1}^K p(x_i | z_i = k, \theta) p(z_i = k | \theta)$$

Gaussian for class k *probability that data point belongs to mixture 'k'!*

k-means: hard-EM
with Gaussian
with $\sum_k = I$

- usually, mixture models are fit with EM

- if one class is more complicated, can use mixture of a Gaussian to get rid of unimodal problem