



<http://www.youtube.com/watch?v=nPdP1jBfxzo>

<http://www.getrobo.com/getrobo/2008/10/keepon-is-now-o.html>

Inference in Bayes nets and Naïve Bayes

CS151
David Kauchak
Fall 2010

Some material borrowed from:
Sara Owsley Sood and others

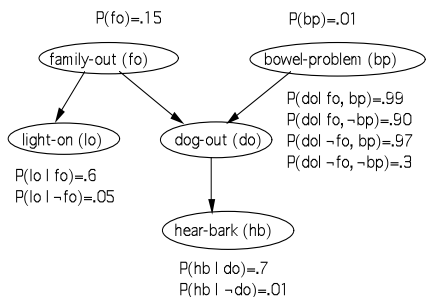
Admin

- Assignment 4 out
- Written 4 and 5
- Grading
- Midterm
- TA office hours

Asking questions about distributions

- We want to be able to ask questions about these probability distributions
- Given n variables, a query splits the variables into three sets:
 - query variable(s)
 - known/evidence variables
 - unknown/hidden variables
- $P(\text{query} \mid \text{evidence})$
 - if we had no hidden variables, we could just multiply all the values in the different CPTs
 - to answer this, we need to sum over the hidden variables!

BN Example



$p(\text{fo} \mid \text{hb}, \text{lo})?$

$p(\text{fo} \mid \text{hb}, \text{lo})$

$$p(\text{fo} \mid \text{hb}, \text{lo}) = \frac{p(\text{fo}, \text{hb}, \text{lo})}{p(\text{hb}, \text{lo})}$$

Evidence: HB, LO
Query: FO
Hidden: BP, DO

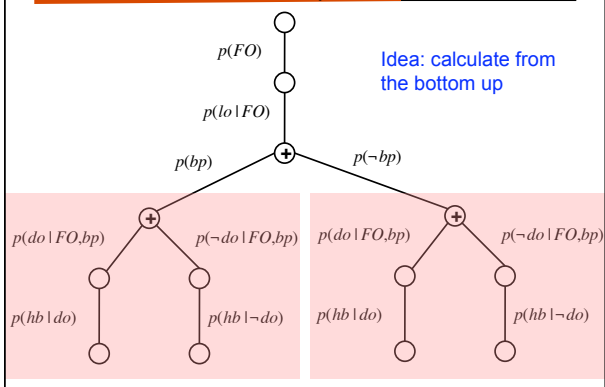
$$p(\text{FO} \mid \text{hb}, \text{lo}) = \alpha p(\text{FO}, \text{hb}, \text{lo})$$

$$= \alpha \sum_{\text{bp}} \sum_{\text{do}} p(\text{FO}, \text{hb}, \text{lo}, \text{bp}, \text{do})$$

$$= \alpha \sum_{\text{bp}} \sum_{\text{do}} p(\text{FO}) p(\text{bp}) p(\text{lo} \mid \text{FO}) p(\text{do} \mid \text{FO}, \text{bp}) p(\text{hb} \mid \text{do})$$

$$= \alpha p(\text{FO}) p(\text{lo} \mid \text{FO}) \sum_{\text{bp}} p(\text{bp}) \sum_{\text{do}} p(\text{do} \mid \text{FO}, \text{bp}) p(\text{hb} \mid \text{do})$$

$$p(\text{FO} \mid \text{hb}, \text{lo}) = \alpha p(\text{FO}) p(\text{lo} \mid \text{FO}) \sum_{\text{bp}} p(\text{bp}) \sum_{\text{do}} p(\text{do} \mid \text{FO}, \text{bp}) p(\text{hb} \mid \text{do})$$



Variable elimination

- Avoids repeated computation
- Break the calculation into *factors*
 - each factor involves some (or all) of the variables
 - factors represent the values for the possible combinations of the variables
 - Initially, these values come straight from the conditional probability tables

$$p(\text{fo}) p(\text{lo} \mid \text{fo}) \sum_{\text{bp}} p(\text{bp}) \sum_{\text{do}} p(\text{do} \mid \text{fo}, \text{bp}) p(\text{hb} \mid \text{do})$$

$f_1(\text{fo})$ $f_2(\text{lo}, \text{fo})$ $f_3(\text{bp})$ $f_4(\text{do}, \text{fo}, \text{bp})$ $f_5(\text{hb}, \text{do})$
 $[0.15]$ $[0.6]$ $\begin{pmatrix} \text{bp} \\ \text{T} & 0.01 \\ \text{F} & 0.99 \end{pmatrix}$ $\begin{pmatrix} \text{do} & \text{fo} & \text{bp} \\ \text{T} & \text{T} & \text{T} & 0.99 \\ \text{T} & \text{T} & \text{F} & 0.90 \\ \text{F} & \text{T} & \text{T} & 0.01 \\ \text{F} & \text{T} & \text{F} & 0.10 \end{pmatrix}$ $\begin{pmatrix} \text{hb} & \text{do} \\ \text{T} & \text{T} & 0.7 \\ \text{T} & \text{F} & 0.01 \end{pmatrix}$

Variable elimination

- What is the size of the factor's matrix dependent on?

- the number of hidden variables, m
- 2^m
- depending on how you treat query variables, they might also factor in here

$$\underbrace{p(fo)}_{f_1(fo)} \underbrace{p(lo|fo)}_{f_2(lo, fo)} \sum_{bp} \underbrace{p(bp)}_{f_3(bp)} \sum_{do} \underbrace{p(do|fo,bp)}_{f_4(do, fo, bp)} \underbrace{p(hb|do)}_{f_5(hb, do)}$$

$$\begin{bmatrix} 0.15 \\ 0.6 \end{bmatrix} \begin{bmatrix} bp \\ T & 0.01 \\ F & 0.99 \end{bmatrix} \begin{bmatrix} do & fo & bp \\ T & T & T & 0.99 \\ T & T & F & 0.90 \\ F & T & T & 0.01 \\ F & T & F & 0.10 \end{bmatrix} \begin{bmatrix} hb & do \\ T & T & 0.7 \\ T & F & 0.01 \end{bmatrix}$$

Variable elimination

$$f_1(fo)f_2(lo, fo) \sum_{bp} f_3(bp) \sum_{do} f_4(do, fo, bp) f_5(hb, do)$$

$$\begin{bmatrix} 0.15 \\ 0.6 \end{bmatrix} \begin{bmatrix} bp \\ T & 0.01 \\ F & 0.99 \end{bmatrix} \begin{bmatrix} do & fo & bp \\ T & T & T & 0.99 \\ T & T & F & 0.90 \\ F & T & T & 0.01 \\ F & T & F & 0.10 \end{bmatrix} \begin{bmatrix} hb & do \\ T & T & 0.7 \\ T & F & 0.01 \end{bmatrix}$$

- Solve this from right to left using two operations:
 - pointwise product of factors
 - summing out a variable

Pointwise product

$$f_1(x_1, \dots, x_n, y_1, \dots, y_m) f_2(y_1, \dots, y_n, z_1, \dots, z_p) = f_3(x_1, \dots, x_n, y_1, \dots, y_n, z_1, \dots, z_p)$$

- When we take the product of two factors, we have three sets of variables
 - x_1, \dots, x_n : those unique to f_1
 - z_1, \dots, z_p : those unique to f_2
 - y_1, \dots, y_n : those shared between the two
- The result is a *new* factor over the union of the variables

Pointwise product

$$f_4(do, fo, bp) f_5(hb, do) = f_6(do, fo, bp, hb)$$

$$\begin{bmatrix} do & fo & bp \\ T & T & T & 0.99 \\ T & T & F & 0.90 \\ F & T & T & 0.01 \\ F & T & F & 0.10 \end{bmatrix} \begin{bmatrix} hb & do \\ T & T & 0.7 \\ T & F & 0.01 \end{bmatrix} = \begin{bmatrix} do & fo & bp & hb \\ T & T & T & T \\ T & T & F & T \\ F & T & T & T \\ F & T & F & T \end{bmatrix}$$

Pointwise product

$$f_4(do, fo, bp)f_5(hb, do) = f_6(do, fo, bp, hp)$$

$\begin{pmatrix} \text{do fo bp} \\ T T T & 0.99 \\ T T F & 0.90 \\ F T T & 0.01 \\ F T F & 0.10 \end{pmatrix}$	$\begin{pmatrix} \text{hb do} \\ T T & 0.7 \\ T F & 0.01 \end{pmatrix}$	$\begin{pmatrix} \text{do fo bp hb} \\ T T T T & 0.7 \cdot 0.99 \\ T T F T & \\ F T T T & \\ F T F T & \end{pmatrix}$
-----------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------

Pointwise product

$$f_4(do, fo, bp)f_5(hb, do) = f_6(do, fo, bp, hp)$$

$\begin{pmatrix} \text{do fo bp} \\ T T T & 0.99 \\ T T F & 0.90 \\ F T T & 0.01 \\ F T F & 0.10 \end{pmatrix}$	$\begin{pmatrix} \text{hb do} \\ T T & 0.7 \\ T F & 0.01 \end{pmatrix}$	$\begin{pmatrix} \text{do fo bp hb} \\ T T T T & 0.7 \cdot 0.99 \\ T T F T & 0.7 \cdot 0.90 \\ F T T T & \\ F T F T & \end{pmatrix}$
-----------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------

Pointwise product

$$f_4(do, fo, bp)f_5(hb, do) = f_6(do, fo, bp, hp)$$

$\begin{pmatrix} \text{do fo bp} \\ T T T & 0.99 \\ T T F & 0.90 \\ F T T & 0.01 \\ F T F & 0.10 \end{pmatrix}$	$\begin{pmatrix} \text{hb do} \\ T T & 0.7 \\ T F & 0.01 \end{pmatrix}$	$\begin{pmatrix} \text{do fo bp hb} \\ T T T T & 0.7 \cdot 0.99 \\ T T F T & 0.7 \cdot 0.90 \\ F T T T & 0.01 \cdot 0.01 \\ F T F T & 0.01 \cdot 0.10 \end{pmatrix}$
-----------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------

In this case the size of the factor didn't increase, but in general, it can

Pointwise Product

A	B	$f_1(A,B)$	B	C	$f_2(B,C)$	A	B	C	$f_3(A,B,C)$
T	T	.3	T	T	.2	T	T	T	0.06
T	F	.7	T	F	.8	T	T	F	0.24
F	T	.9	F	T	.6	T	F	T	0.42
F	F	.1	F	F	.4	T	F	F	0.28
						F	T	T	0.18
						F	T	F	0.72
						F	F	T	0.06
						F	F	F	0.04

Summing out a variable

$$f_1(fo)f_2(lo,fo)\sum_{bp}f_3(bp)\sum_{do}f_6(do,fo,bp,hb)$$

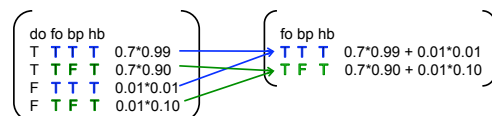
$$\begin{bmatrix} 0.15 \\ 0.6 \end{bmatrix} \begin{bmatrix} bp \\ T \ 0.01 \\ F \ 0.99 \end{bmatrix} \begin{bmatrix} do \ fo \ bp \ hb \\ T \ T \ T \ T \ 0.7*0.99 \\ T \ T \ F \ T \ 0.7*0.90 \\ F \ T \ T \ T \ 0.01*0.01 \\ F \ T \ F \ T \ 0.01*0.10 \end{bmatrix}$$

- Produces a new factor with one less variable
- Reduces the size of the table by a factor of the number of possible values for the variable (for binary 2)

Summing out a variable

$$\sum_{do} f_6(do, fo, bp, hb) = f_6(do, fo, bp, hp) + f_6(-do, fo, bp, hp)$$

$$= f_7(fo, bp, hp)$$



How do we sum out a variable?

Variable Elimination

$$f_1(fo)f_2(lo,fo)\sum_{bp}f_3(bp)f_7(fo,bp,hb)$$

$$\begin{bmatrix} 0.15 \\ 0.6 \end{bmatrix} \begin{bmatrix} bp \\ T \ 0.01 \\ F \ 0.99 \end{bmatrix} \begin{bmatrix} fo \ bp \ hb \\ T \ T \ T \ 0.7*0.99 + 0.01*0.01 \\ T \ T \ F \ T \ 0.7*0.90 + 0.01*0.10 \\ F \ T \ T \ T \ 0.01*0.01 \\ F \ T \ F \ T \ 0.01*0.10 \end{bmatrix}$$

$$f_1(fo)f_2(lo,fo)\sum_{bp}f_8(fo,bp,hb) \quad \text{product}$$

$$f_1(fo)f_2(lo,fo)f_9(fo,hb) \quad \text{sum}$$

⋮

Variable ordering

$$f_1(fo)f_2(lo,fo)\sum_{bp}f_3(bp)\sum_{do}f_4(do,fo,bp)f_5(hb,do)$$

$$=$$

$$f_1(fo)f_2(lo,fo)\sum_{do}f_5(hb,do)\sum_{bp}f_3(bp)f_4(do,fo,bp)$$

- The complexity depends on which order we sum out the variables

Variable ordering

$$\sum_a \sum_b f_1(a,b) f_2(b,c) \quad \text{vs} \quad \sum_b f_2(b,c) \sum_a f_1(a,b)$$

$$\underbrace{\sum_a \sum_b f_3(a,b,c)}$$

A factor containing 3 variables

$$\underbrace{\sum_b f_2(b,c) f_3(b)}$$

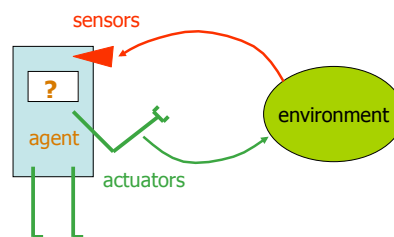
A factor containing only 2 variables

Runtime

- In general, the run-time of the variable elimination algorithm is dependent on the largest factor created
- Figuring out the optimal variable ordering is intractable
- Some heuristics have been used
 - pick the merger greedily

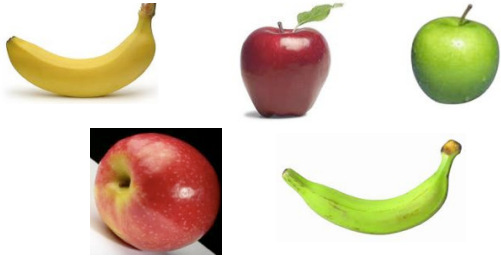
Learning from Data

Learning



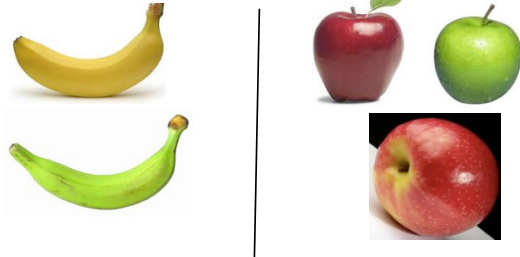
As an agent interacts with the world, it should learn about its environment

Lots of different learning problems



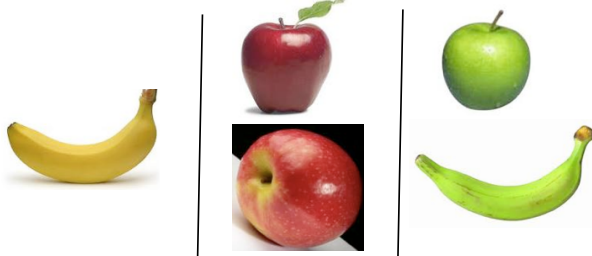
Unsupervised learning: put these into groups

Lots of different learning problems



Unsupervised learning: put these into groups

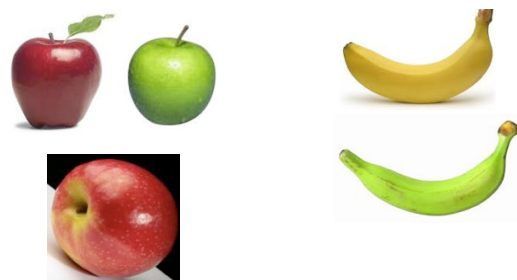
Lots of different learning problems



No explicit labels/categories specified

Unsupervised learning: put these into groups

Lots of learning problems



APPLES

BANANAS

Supervised learning: given labeled data

Lots of learning problems

- Given labeled examples, learn to label unlabeled examples



APPLE or BANANA?

Supervised learning: learn to classify unlabeled

Lots of learning problems

- Many others
 - semi-supervised learning: some labeled data and some unlabeled data
 - active learning: unlabeled data, but we can pick some examples to be labeled
 - reinforcement learning: maximize a *cumulative* reward. Learn to drive a car, reward = not crashing
- and variations
 - online vs. offline learning: do we have access to all of the data or do we have to learn as we go
 - classification vs. regression: are we predicting between a finite set or are we predicting a score/value

Supervised classification: training

Labeled data

Data Label



0



0



1



1



0

train a predictive model



Supervised learning: testing/classifying

Unlabeled data



predict the label

labels

1

0

0

0

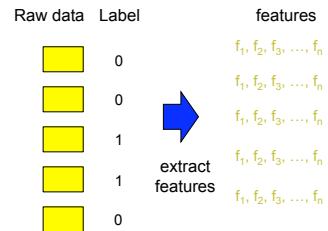
1

0

Some example

- image classification
 - does the image contain a person? apple? banana?
- text classification
 - is this a good/bad review?
 - is this article about sports or politics?
 - is this e-mail spam?
- character recognition
 - is this set of scribbles an 'a', 'b', 'c', ...
- credit card transactions
 - fraud or not?
- audio classification
 - hit or not?
 - jazz, pop, blues, rap, ...
- Tons of problems!!!

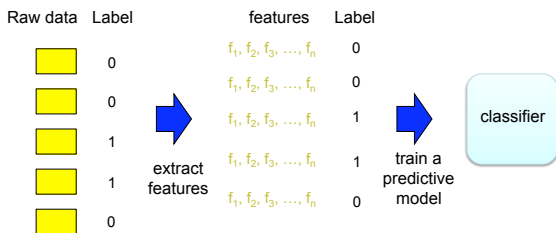
Features



- We're given "raw data", e.g. text documents, images, audio, ...
- Need to extract "features" from these (or to think of it another way, we somehow need to represent these things)
- What might be features for: text, images, audio?

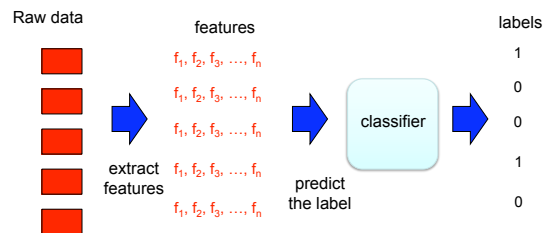
Feature based classification

Training or learning phase



Feature based classification

Testing or classification phase



Bayesian Classification

We represent a data item based on the features:

$$D = \langle f_1, f_2, \dots, f_n \rangle$$

Training

$$\begin{aligned} \text{a: } p(a|D) &= p(a|f_1, f_2, \dots, f_n) \\ \text{b: } p(b|D) &= p(b|f_1, f_2, \dots, f_n) \end{aligned} \rightarrow P(\text{Label}|f_1, f_2, \dots, f_n)$$

For each label/class, **learn** a probability distribution based on the features

Bayesian Classification

We represent a data item based on the features:

$$D = \langle f_1, f_2, \dots, f_n \rangle$$

Classifying

$$\text{label} = \underset{l \in \text{Labels}}{\text{argmax}} P(l|f_1, f_2, \dots, f_n)$$

Given an *new* example, classify it as the label with the largest conditional probability

Bayes rule for classification

$$P(\text{Label} | \text{Data}) = \frac{P(D|C)P(C)}{P(D)}$$

conditional (posterior) probability
prior probability

Bayesian Classifiers

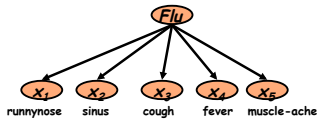
$$\text{label} = \underset{l \in \text{Labels}}{\text{argmax}} P(l|f_1, f_2, \dots, f_n) \leftarrow \text{different distributions for different labels}$$

$$= \underset{l \in \text{Labels}}{\text{argmax}} \frac{P(f_1, f_2, \dots, f_n | l)P(l)}{P(f_1, f_2, \dots, f_n)} \quad \text{Bayes rule}$$

$$= \underset{l \in \text{Labels}}{\text{argmax}} P(f_1, f_2, \dots, f_n | l)P(l)$$

two models to learn for each label/class

The Naive Bayes Classifier



Conditional Independence Assumption: features are independent of each other given the class:

$$P(x_1, \dots, x_n | l) = P(x_1 | l)P(x_2 | l) \dots P(x_n | l)$$

$$label = \operatorname{argmax}_{l \in Labels} P(f_1 | l)P(f_2 | l) \dots p(f_n | l)P(l)$$

Estimating parameters

- I flip a coin 1000 times, how would you estimate the probability of heads?
- I roll a 6-sided die 1000 times, how you estimate the probability of getting a '6'?

For us:

$$label = \operatorname{argmax}_{l \in Labels} P(f_1 | l)P(f_2 | l) \dots p(f_n | l)P(l)$$

Ideas?

Maximum likelihood estimates

$$\hat{P}(l) = \frac{N(l)}{N} \quad \begin{array}{l} \text{number of items with label} \\ \text{total number of items} \end{array}$$

$$\hat{P}(f_i | l) = \frac{N(f_i, l)}{N(l)} \quad \begin{array}{l} \text{number of items with the label with feature} \\ \text{number of items with label} \end{array}$$

Any problems with this approach?

Problem with Max Likelihood

- What if we have seen no training cases where patient had no flu and muscle aches?

$$\hat{P}(f_{mu} | n_f) = \frac{N(f_{mu}, n_f)}{N(n_f)} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$label = \operatorname{argmax}_{l \in Labels} \hat{P}(l) \prod_i \hat{P}(f_i | l)$$

Smoothing to Avoid Overfitting

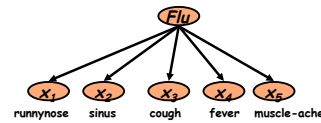
Make every event a little probable...

$$\hat{P}(f_i | l) = \frac{N(f_i, l) + \lambda}{N(l) + k\lambda}$$

of features

Unseen features

- Note that this is different from coming in with a feature we've never seen before (in any of the classes)
 - For example, "bloating"



Naïve Bayes Text Classification

- Features: word occurring in a document (though others could be used...)

$$label = \underset{l \in Labels}{\operatorname{argmax}} P(word_1 | l) P(word_2 | l) \dots p(word_n | l) P(l)$$

- Does the Naïve Bayes assumption hold?
 - Are word occurrences independent given the label?
- We'll look at a few application for this homework
 - sentiment analysis: positive vs. negative reviews
 - category classification

Classification evaluation?

- Accuracy
 - num correct / total
- Class specific measures
 - Precision
 - num correct with class A / num predicted class A
 - Recall
 - num correct with class A / num with class A
 - F1-measure
 - $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$
- Why have these class specific measures?

WebKB Experiment (1998)

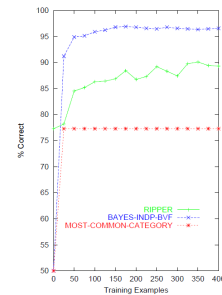
- Classify webpages from CS departments into:
 - student, faculty, course, project
- Train on ~5,000 hand-labeled web pages
 - Cornell, Washington, U.Texas, Wisconsin
- Crawl and classify a new site (CMU)

- Results:



	Student	Faculty	Person	Project	Course	Department
Extracted	180	66	246	99	28	1
Correct	130	28	194	72	25	1
Accuracy:	72%	42%	79%	73%	89%	100%

Naive Bayes on spam email



<http://www.cncb.cmu.edu/~jp/research/email.paper.pdf>

SpamAssassin

- Naive Bayes has found a home in spam filtering
 - Paul Graham's *A Plan for Spam*
 - A mutant with more mutant offspring...
 - Naive Bayes-like classifier with weird parameter estimation
 - Widely used in spam filters
 - But also many other things: black hole lists, etc.
- Many email topic filters also use NB classifiers

NB: The good and the bad

- Good
 - Easy to understand
 - Fast to train
 - Reasonable performance
- Bad
 - We can do better
 - Independence assumptions are rarely true
 - Smoothing is challenging
 - Feature selection is usually required