5.1 a) why prefer an unbiased estimator to a biased estimator?
the mean squared error is minimized where
$$\mu_E = \mu$$

b) does an unbiased statistic ensure its a good estimator?
it is one of two criteria, the other is for it
to be "stable" where the variation of $E$ is small
$\sigma_E^2 << 0$. If another random sample is taken
we dont want $\bar{x}$ to be random

c) when biased might be a better choice?
if the bias can reduce the variance
such that the mean square error decreases

5.2 normal distribution, std dev = 0.5 mm
ses lengths: 75.3, 76.0, 75.0, 77.0, 75.4, 76.3, 77.0, 74.9, 76.5, 75.8
a) what is the parameter of interest?
The mean length of changed parts
b) 759.2/10
75.92
c) .99 confidence
lower    upper    $\mu E = .407$
(75.51, 76.33)

5.3 why is it necessary to test a classifier with independent data
The second alg is better because the 5% error in
training data allows for learning & therefore has
less error for test data. A classifier needs to
be tested with independent data to increase the
generalizeing capability of the model for unseen data.
$e = .632\, e_{test} + .568\, e_{resub}$
$e_1 = 0 + .368(.2) = .0736$        $min(e_1, e_2) = e_2$
$e_2 = .632(.05) + .368(.1) = .0684$  therefore 2nd alg is better

Predicted class

|  | | a | b | c | d |
|---|---|---|---|---|---|
| S.4 | a | 15 | 2 | 3 | 5 |
| actual | b | 1 | 18 | 6 | 2 |
| class | c | 3 | 5 | 9 | 0 |
| | d | 2 | 2 | 7 | 12 |

$R_s = 1 - R_e$

$$\frac{15 + 18 + 9 + 12}{1 + 3 + 2 + 2 + 5 + 2 + 7 + 3 + 6 + 5 + 2 + 0 + 54}$$

$$\frac{54}{38 + 54} = .587 = R_s$$

$$R_e = .413$$

S.5   # cases = 609

   Cancer = 84

   not cancer = 525

|  | C | WC |
|---|---|---|
| C | TP | FP |
| NC | FN | TN |

total # of positive = TP + FN

Success = 75%

Error = 25%

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{FP + TN}$$

FP = 525 · .25 = 131

TN = 525 · .75 = 354

$$Sens = \frac{75}{75 + 9} = .892$$

$$Spec = \frac{354}{525} = .674$$

7.5    $\dfrac{423}{21} = 20.14 = \mu$

median $= 16$
outliers $= 69$ & $55$

    $\dfrac{299}{19} = 15.74$    median $\neq 6$

A will have a higher avg than b because
the two outliers are in the upper bound
of samples which increases the mean
the median is the same as outliers
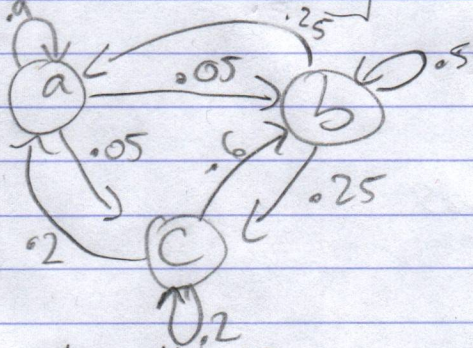do not out vergh the # of middle
#s

7.6 Class code 2 is better because more
    bits offers a more accurate representation
    of teamming distance
    B would be correct class as they have
least # of diferences

a)

$$
\begin{array}{c}
\quad\quad a \quad\quad b \quad\quad c \\
\begin{array}{c} a \\ b \\ c \end{array}
\left[\begin{array}{ccc}
.9 & .05 & .05 \\
.25 & .5 & .25 \\
.2 & .6 & .2
\end{array}\right] = P
\end{array}
$$



b)

$$
\begin{array}{ccc}
L & L1 & L2
\end{array}
$$
$$S_0 = [.8 \quad .1 \quad .1]$$

$$
\underset{1 \times 3}{\quad} \quad\quad\quad \underset{3 \times 3}{\quad}
$$

$$
S_1 = S_0 P = [.8 \quad .1 \quad .1]
\left[\begin{array}{ccc}
.9 & .05 & .05 \\
.25 & .5 & .25 \\
.2 & .6 & .2
\end{array}\right] =
$$

$$
\begin{array}{ccc}
2 & L1 & L2
\end{array}
$$
$$S_1 = [.765 \quad .15 \quad .085]$$

$$
\begin{array}{ccc}
c & c1 & L2
\end{array}
$$
$$S_1 P = S_2 = [.743 \quad .16425 \quad .09275]$$

$$S_6 = S_0 P^6 = [.7078 \quad .1900 \quad .1025]$$

$$S_0 P^{99} = [.6993 \quad .1958 \quad .1049]$$

$$K = 15 \quad SP = S$$

$$
\begin{array}{c}
\quad\quad\quad\quad S \quad\quad\quad\quad\quad\quad\quad\quad\quad P \\
S_{15} = [.6995 \quad .1956 \quad .1048]
\left[\begin{array}{ccc}
.9 & .05 & .05 \\
.25 & .5 & .25 \\
.2 & .6 & .2
\end{array}\right] =
\end{array}
$$

$$SP = S = [.6994 \quad .1957 \quad .1049]$$