

Data Mining, CSCI 347, Fall 2019
Evaluation – Counting Costs, Nov. 13

1. A confusion matrix for a model is given. Tell the accuracy of the model.

		Predicted class			<i>total</i>
		<i>a</i>	<i>b</i>	<i>c</i>	
Actual class	<i>a</i>	88	10	2	100
	<i>b</i>	14	40	6	60
	<i>c</i>	18	10	12	40
<i>total</i>		120	60	20	

2. The same confusion matrix is given below on the left. On the right is a confusion matrix for a random predictor. Calculate the Kappa statistic to measure the relative improvement of the model over a random predictor.

Model Predictions

Random Predictor

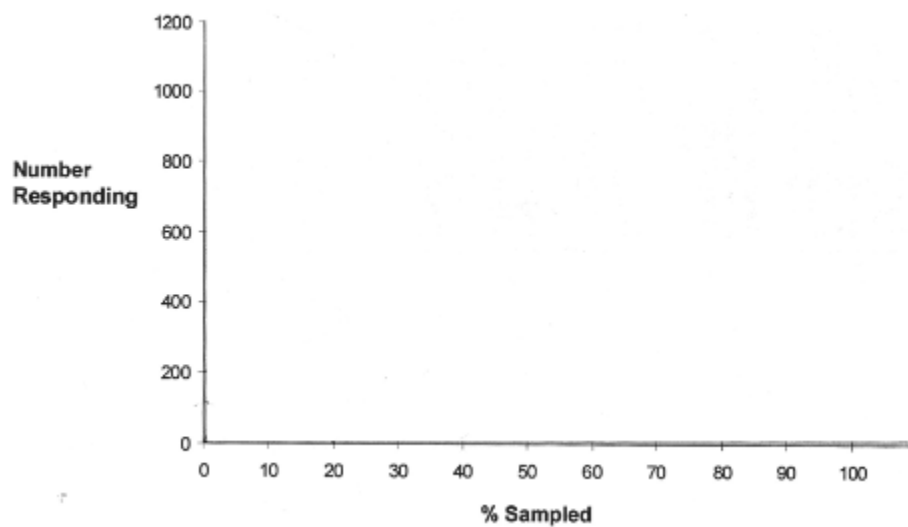
		Predicted class								Predicted class			
		<i>a</i>	<i>b</i>	<i>c</i>	<i>total</i>					<i>a</i>	<i>b</i>	<i>c</i>	<i>total</i>
Actual class	<i>a</i>	88	10	2	100			<i>a</i>	60	30	10	100	
	<i>b</i>	14	40	6	60		Actual class	<i>b</i>	36	18	6	60	
	<i>c</i>	18	10	12	40		<i>c</i>	24	12	4	40		
<i>total</i>		120	60	20			<i>total</i>	120	60	20			

3. Example: Say that we are going to do a direct mailing and we have a million addresses in our database. In general we know that 0.1% (0.001) household will respond to our mailing. Say that using data mining, we have an algorithm that identifies a subset of 100,000 of the most promising addresses, where these households are likely to respond at a rate of 0.4%. Another model identifies a subset of 400,000 household where 0.2% are likely to respond.

Create a lift chart to visualize this information.

x axis is the percent sampled

y axis is number of true positives



4. Say when predicting the recurrence of breast cancer, 70% of the women that survived breast cancer had not had recurrence within 5 years. The other 30% had. Call the recurrence of breast cancer a positive. For parts a and b, say that there are 100 women.

a. Calculate the F1 score of a model that always predicts that there will be no recurrence.

b. Calculate the F1 score if a model always predicts that there will be a recurrence.

c. CART (Classification And Regression Trees) for 300 women, has the confusion matrix:

	Predicted		
	Recurrence	No Recurrence	
Actual	77	13	
	No Recurrence	39	171

Calculate the F1 score.