

Exam 1 SG

Friday, February 23, 2018 8:54 PM

- If you are ever working in a data science role, you have to learn the CRISP-DM methodology, which is a process model. We will tie different aspects of the CRISP-DM cycle into the projects of the course.
 - o The business understanding is where you are figuring out what you are trying to solve for in order to make scientific progress.
 - o The data understanding is when you try to find and comprehend the data you have to work with.
 - o The data preparation is when you modify the data you have collected to fit the record structure of your modeling process
 - o The modeling process is where you determine the function that describes new knowledge based on previous attributes
 - o The evaluation process is when you evaluate the accuracy of your model, as well as the quality of the problem presenters
 - o The deployment process is how you utilize your new model for exterior gain.
- Classifying takes in a set of data, and splits it into training and test data. It assigns a class attribute to each row based on the attributes in the row. It takes training data that has preassigned classes, and builds a model that can predict the class of the test data. The test data is then used to determine the goodness of the model. The end goal is to be able to predict the class of unseen records with a certain degree of confidence.
- Clustering finds data points that are similar in one another according to their attributes. These attributes could be things like Euclidian distance or other things. The data points that are similar are all assigned to the same cluster. The goal is to find all the clusters for a data set, and be able to assign new attribute vectors to existing clusters.
- Association rules take a set of records with a number of items from a given collection, and produce rules that will predict if another item will appear in the collection based on previous items.
- Data preprocessing takes un/semi structured data and turns it into structured data. This means that raw data from databases and text documents turn into things like R dataframes and tables.
- There are types of attributes:
 - o Nominal, like ID numbers or eye color. These are things that just represent identifying names, and are used to distinguish one attribute type from another.
 - o Ordinal, which are rankings. These allow you to order attributes in a specific way, and compare them.
 - o Intervals, which are things where intervals between attribute values make sense, like date ranges and distances between attribute values.
 - o Ratios, which are things where both intervals and ratios make sense to look at. You would be able to multiply and divide ratios.

which is a data mining

course.

or financial gain or

work with.

structure used in the

ge from the data using

gain it delivers to your

route based on a subset of

el based on the training

e able to classify previously

re similarity measures can

ed to one cluster. We want

usters.

roduce dependency rules

ns things like relational

ames for data. Allows you

ompare items.

ates. You can look at times

uldn't say that march 2nd is

twice as large as March 1st, but you would say that 2 minutes is twice as long as 1 minute. Measurement variables are intervals, but with the added condition that 0 of the measurement variable means none of that variable.

- Attributes can be discrete or continuous. Discrete attributes are finite or countably infinite, and continuous attributes are ordinal. Continuous variables have real numbers as attribute values, and are typically interval or ratio variables.
- The dot product of two vectors is just the sum of the product of their corresponding entry values. The dot product of $\mathbf{a} = [a_1 \ a_2 \ \dots \ a_m]^T$ and $\mathbf{b} = [b_1 \ b_2 \ \dots \ b_m]^T$ is just $a_1b_1 + a_2b_2 + \dots + a_mb_m$.

Length

The *Euclidean norm* or *length* of a vector $\mathbf{a} \in \mathbb{R}^m$ is defined as

$$\|\mathbf{a}\| = \sqrt{\mathbf{a}^T \mathbf{a}} = \sqrt{a_1^2 + a_2^2 + \dots + a_m^2} = \sqrt{\sum_{i=1}^m a_i^2}$$

Distance

From the Euclidean norm we can define the *Euclidean distance* between \mathbf{a} and \mathbf{b} , as follows

$$\delta(\mathbf{a}, \mathbf{b}) = \|\mathbf{a} - \mathbf{b}\| = \sqrt{(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})} = \sqrt{\sum_{i=1}^m (a_i - b_i)^2} \quad (1.1)$$

Angle

The cosine of the smallest angle between vectors \mathbf{a} and \mathbf{b} , also called the *cosine similarity*, is given as

$$\cos\theta = \frac{\mathbf{a}^T \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} = \left(\frac{\mathbf{a}}{\|\mathbf{a}\|} \right)^T \left(\frac{\mathbf{b}}{\|\mathbf{b}\|} \right) \quad (1.3)$$

Thus, the cosine of the angle between \mathbf{a} and \mathbf{b} is given as the dot product of the unit vectors $\frac{\mathbf{a}}{\|\mathbf{a}\|}$ and $\frac{\mathbf{b}}{\|\mathbf{b}\|}$.

- A percentile, given an ordinal or continuous attribute x and a number p between 0 and 100, the x_p of x such that $p\%$ of the observed values of x are less than x_p .

- ❖ The **variance** is the most common measure of the spread of a set of points.

$$\text{variance}(x) = s_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

- ❖ The **standard deviation** is the square root of the variance.
 - ❖ However, this is also sensitive to outliers, so that other measures are often used.

$$\text{AAD}(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$$

te. Jason says that ratio
e indicated that there is

nd are typically nominal or
or ratio typed.

ues. So $[a_1, a_2] \text{ dot } [b_1, b_2]$

the pth percentile is a value

$$\text{MAD}(x) = \text{median} \left(\{|x_1 - \bar{x}|, \dots, |x_m - \bar{x}|\} \right)$$

$$\text{interquartile range}(x) = x_{75\%} - x_{25\%}$$

- Covariance looks at every attribute value in a dataframe of two attributes. It is the sum of the corresponding attribute value divided by the number of rows -1. Think of it like scrolling through multiplying the two attribute values together into a new attribute, summing that new attribute dividing by the number of rows in that vector -1.

$$s_{j,k} = \text{cov}(x_j, x_k) = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

- The covariance matrix is the matrix comprised of every possible covariance in your dataset. So of the matrix is your first attribute covaried with all the other attributes, the second row is the second attribute, and so on and so forth. The correlation between two attributes is their entry (their covariance) divided by the square root of the products of their individual variances:

$$r_{j,k} = \frac{s_{j,k}}{\sqrt{s_{j,j}s_{k,k}}}.$$

- Correlation measures the linear relationship between two continuous attributes.

Bayes Theorem

$$P(Y|X) = \frac{P(X|Y)}{P(X)} P(Y)$$

Posterior Probability = Support x Prior Probability

1. Start with some prior “belief” ($P(Y)$) and observe data X .
2. $P(X|Y)$ is the likelihood of observing X given our belief Y . You want to normalize this by the likelihood of observing X (regardless of Y). The ratio of these two numbers is the support.
3. Update your belief in Y with this new information ($P(Y|X)$, the posterior probability).

► **Key Idea:** If the support is greater than 1, then the observed data X will increase your belief in Y

products of every
row through the dataframe,
sum the vector up, and then

do everything in the first row
the same thing but with the
y in the covariance matrix

- $P(X) = P(X|Y)P(Y) + P(X|!Y)P(!Y)$.
- So we adjust our prior/initial belief in Y by first assuming we did get Y, then testing X based on to get $P(X)$, we look at the opposite, where we look at the probability of not getting Y and what on not getting Y. Then we use the formula to get our updated belief in X given what we know
- **Proximity** refers to either similarity or dissimilarity between points/attribute values. Distance dissimilarity. To scale dissimilarities, use this formula:
 - o $(\text{original dissim} - \text{minimum dissim value}) / (\text{max dissim value} - \text{min dissim value})$, where values from your original range of dissimilarities.
- The Euclidian distance is:

o

$$\text{dist} = \left(\sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

- o Where $r = 2$. You iterate through the attribute values, adding them together after subtracting them. Then, you square root the result.
- The simple matching coefficient is used only for binary attributes. It is denoted as the number of negatives divided by the number of all outcomes. Jaccard is an offshoot, and only calculates the matches, divided by the number of all attributes except for true negative matches. It is more appropriate for vectors with a lot of matching true negatives:

$$\text{SMC} = \text{number of matches} / \text{number of attributes}$$

$$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$

$$J = \text{number of 11 matches} / \text{number of not-both-zero attributes}$$

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

- Cosine similarity treats the two attribute rows in the dataframe as vectors. It is the dot product divided by the product of their norms:

- $\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|$

- Remember that with proximity, different attributes have different scales so you should scale your computing distance.

n that to get $P(X|Y)$. Then,
at $P(X|Y)$ would be based
about Y.

is a common term for

the max and mins are the

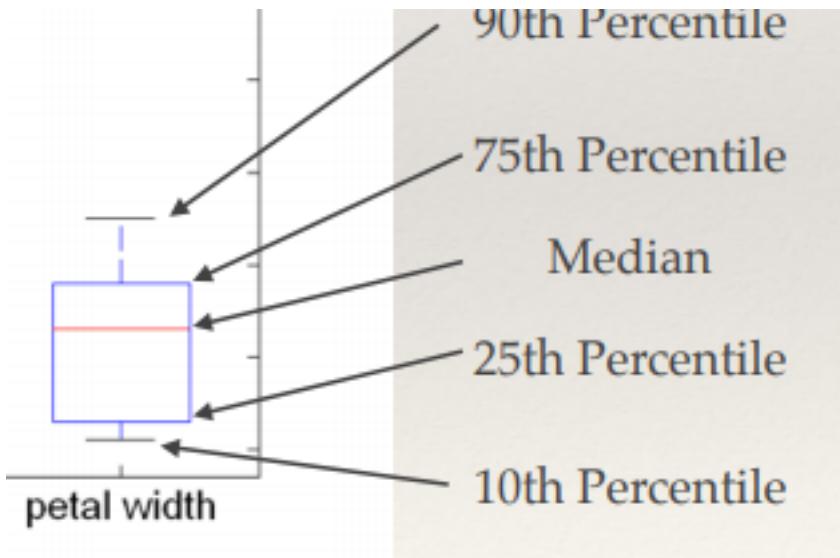
acting them, squaring

or of true positives and true
the number of true positive
appropriate for attribute

utes values

ct of the two vectors

your data before



- **Noise** in data refers to modification of original values, it can occur to measurement error or numerous other things. It creates flawed biased data.
- Outliers and missing values are also big problems in data mining.

Data Preprocessing

- This is a technique used to make data more suitable for data mining. You select data objects and/or create/change attributes so that mining is easier. There are many types:
 - o Aggregation combines two or more attributes into a single attribute, so as to reduce dimensions and create more stable data with less variability. You can lose interesting data this way too.
 - o Sampling is done by only looking at one or a couple of portions from an entire dataset. You can do this with replacement, so as to repeat draws or not from your data. If the sample is representative of the entire dataset, it will work almost as well as if you were using the entire thing.
- The curse of dimensionality is when you have higher dimensions of attributes in your dataframes. As the number of dimensions increases, the volume of the space increases exponentially, which makes the data sparse. Distance and density become less meaningful.
- Dimensionality Reduction methods are done to reduce the number of attributes in one row or column by creating new attributes as combinations of your old attributes.

Principle Component Analysis

- An eigenvector is a vector associated with a transformation matrix that does not change when multiplied by the matrix, just magnitude. An $N \times N$ matrix has N eigenvectors. You can scale it by a scalar and multiply it by the transformation matrix and it will remain the same, just larger in magnitude. Also, eigenvectors are orthogonal to each other. If you find an eigenvector with a magnitude of 1. So if you find a bigger one, you just scale down.
- Eigenvalues are associated with eigenvectors. The amount the eigenvector is scaled up by when multiplied by the transformation matrix is the eigenvalue. So if the eigenvector was 1, 1 before and becomes 4, 4 after, the eigenvalue is 4.
- To find the eigenvectors of a matrix, create a candidate set of likely vectors, multiply them by the matrix, and see which entries all scale up by the same amount. That amount is your eigenvalue, and the vector is your eigenvector!

uman error, among

nd attributes for analysis,

mensionality, change scale,

You can sample with or
esentative of your entire

ne, the available data

ful.

f your dataframe by

change orientation when
before multiplying with
ctors are usually expressed

en multiplying by the
4, then the eigenvalue is 4.
your square matrix, and
ector pre scaling is an

- PCA is a strategy to compress data from a number of dimensions to a smaller amount of dimensions. It still contain a lot of the interesting variation in the data, but are easier to work with as less dimensions.
 - o Step 1 is to get data.
 - o Step 2 is to subtract the average of each attribute column from every element in that column.
 - o Step 3 is to Calculate the covariance matrix of your dataset
 - o Step 4 is to Find the eigenvalues and eigenvectors of the covariance matrix, and make sure that the eigenvalues/vectors, IE, their magnitude is 1.
 - o In step 5, the eigenvector with the highest eigenvalue is the principle component of the dataset. This is the vector that most closely points down the middle of both attributes on a scatterplot of your data. You can then order the eigenvectors by the value of their eigenvalues, and construct a matrix of the eigenvectors (the feature vector) you want to keep. If an eigenvector has a low eigenvalue, you can eliminate it from the matrix without losing a lot of information about your original dataset.
 - o Step 6 involves deriving the new dataset. You transpose your dataset and the feature vector, and then multiply everything together. This expresses your dataset in the terms of your feature vector, which is now the first dimension of your feature vector. So if you have a 1 eigenvector feature vector, you reduced the dimension.

Classification/Model Evaluation

- The input data for a classification task is a tuple (x, y) , where x is a row of attributes, and y is the class label.
- Classification is the task of learning a function f that maps each attribute set x to the class y .
- Class labels are pretty much only binary or nominal.
- The generalization capability of a classifier is its ability to correctly map previously unseen data. To do this, it creates a model based on a training set, which is a dataset with a set of attributes identified as training data. This is a set of rows of attributes that must be classified.
- To test the performance of our model, we have to use a confusion matrix, which denotes whether each row was correctly classified or not:

		PREDICTED CLASS	
		Class +	Class -
ACTUAL CLASS	Class +	a (TP)	b (FN)
	Class -	c (FP)	d (TN)

a: TP

b: FN

c: FP

d: TN

nsions. Those dimensions
dimensions = less work.

olumn.

ure they are the unit

dataset. This represents a
our attribute data. You
eigenvectors (feature
om this matrix without

ector, and multiply
th the dimension being the
duce your dataset to one

the class assigned to that

a to its class label. To do
tical to the test set, which

ther a class in the test set

’ (true positive)

N (false negative)

’ (false positive)

		(TP)	(TN)
Class -		c (FP)	d (TN)

d: TN

- True positive is when it correctly classifies something, as well as true negative. False negative errors.

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Error Rate = $(FN + FP) / (TP + FN + FP + TN)$ =

True Positive Rate (TPR) or Sensitivity = $TP / (TP + FN)$ =

True Negative Rate (TNR) or Specificity = $TN / (TN + FP)$ =

False Positive Rate (FPR) = $FP / (FP + TN)$ = $1 - TNR$

False Negative Rate (FNR) = $FN / (FN + TP)$ = $1 - TPR$

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

N (true negative)

and false positive are

$$\frac{TN}{FP + FN}$$

1-Accuracy

$$\frac{(TP + FN)}{(TN + FP)}$$

TNR

1 - TPR

Decision Trees

- Decision tree classifiers have three elements, a root node making an initial decision, internal nodes representing subsequent decisions, and leaf nodes representing ultimate classifications of rows. Root and internal nodes have two or more edges pointing out representing different questions about a single attribute in the row.
- Hunts Algorithm is a way to recursively create a decision tree:
 - o Step 1 is if all rows in a set of rows with class labels are of the same class label, create a leaf node.
 - o Step 2 is if all the rows in a set of rows with class labels are not of the same class label, find a condition to split them up into separate sets of data with minimal difference in the class labels and apply the algorithm to all the new separate sets of data.
- So we need both a good way to know how to split test data into different subsets, and to know where to place new nodes in the tree. First we will look at splitting on attributes.
 - o We can split nominal and ordinal attributes into their own edges, or groups of them into intervals.
 - o We can split continuous data into intervals, like $<$, $>$, or $x < y < z$.
- There are measures we can use to determine the best way to divide rows into different sets based on these splits. These are called Entropy, Gini, and Classification Error. They measure leaf node impurity. The splits from these measurements are usually the best possible split based on an attribute.
- The Gini index is written as:

$$GINI(t) = 1 - \sum_j [p(j | t)]^2$$

- Where $p(j | t)$ is simply the amount of a type of class in one of the resulting subsets after an attribute split: $(\text{class 0} / \# \text{ rows in subset})^2 - \dots - (\text{class } n / \# \text{ rows in subset})^2$. To figure out the quality of a split, you want to minimize GINI:

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where, n_i = number of records at child i ,

n = number of records at node p .

nodes which represent internal nodes can have one row.

leaf node for this class

then find some optimal test labels. Then, go through

now when to stop creating

to their own edges.

based on attribute edges.
lowest value of impurity

attribute split. So it is $1 - \text{using GINI}$, use this

- When splitting continuous attributes, line all of the attributes up into a row sorted ascending. Then, test for \leq and $>$ on all the continuous values that are in between the sorted actual values. Then, test for \leq and $>$ on all the continuous values that are in between the sorted actual values.

Cheat	No	No	No	Yes	Yes	Yes
	Taxable Income					
	60	70	75	85	90	95
	55	65	72	80	87	92
	\leq	$>$	\leq	$>$	\leq	$>$
Yes	0	3	0	3	0	3
No	0	7	1	6	2	5
Gini	0.420	0.400	0.375	0.343	0.417	0.400
	0.300	0.275	0.250	0.225	0.200	0.175

- Another metric you can use to split up attributes and look at how pure they are is entropy. Instead of Gini score, it uses the information gain score to judge a split (you want the lowest value for pure entropy).

$$\text{Entropy}(t) = -\sum_j p(j | t) \log_2 p(j | t)$$

$$GAIN_{split} = Entropy(p) - \left(\sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

- Information gain is defined as the above, Where k is the number of attributes values the split creates, n is the total number of records in the parent node you are splitting up, ni is the number of records in the partitioned child internal node, and Entropy(i) is the entropy of the partitioned child internal node. You want the largest information gain.
- Problems can arise when using info gain, as splitting into a large number of sets with very few records in each lead to less reliable predictions, so we use the gain ratio to adjust information gain by penalizing splits with many partitions.

Then, find discrete values
ues:

No	No	No	No					
100	120	125	220					
	110	122	172	230				
>	\leq	$>$	\leq	$>$	\leq	$>$		
0	3	0	3	0	3	0		
4	4	3	5	2	6	1	7	0
00	0.343	0.375	0.400	0.420				

instead of the GINI split
(entropy):

$$\log_2 p(j | t)$$

(i)

can take on, n is the
partitioned child internal
st gain possible.
records in each. This can
ing large numbers of small

partitions.

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO}$$

$$SplitINFO =$$

- Error is another metric:

$$Error(t) = 1 - \max_i P(i)$$

- To measure the gain in error, just use information gain but with error stuff instead.
- If you have to find the information gain of an entire dataset of data, Entropy(p) is actually the classes of the rows.
- There are two ways to handle tree overfitting when inductively growing a tree. The first is to p some good stopping condition for when you do not need to split up an impure node any further involve if all classes or attribute values are the same, or if the number of rows in the new node specified constant. Another stopping measure is if the Gini split or info gain score is not impro created. This can be problematic.
- The other option is post pruning. In post pruning, we try to simplify the decision tree to avoid two ways, through subtree replacement, where we replace a subtree in the tree with its major raising, where we replace the contents of an entire subtree with its most often used branch. A split is bad is to look at the error before and after splitting, and if the error is higher, you prun node using subtree replacement with the majority class.
- Decision trees can suffer from data fragmentation, as the leaf nodes could be too small to ma significant decision. Also, finding an optimal tree is NP Hard.

Estimating Generalization Error

- It is generally agreed that the complexity of a model has a factor in its tendency to overfit data to estimate the generalization error, which is the error that will occur when classifying unknown ways to do this, using Resubstitution Error (the error of the dataframe used to train the mode error, pessimistic error, and validation set error).
- Optimistic error just assumes that the error for the test set will be the same as the error for th your training data. This is typically not the best way to handle error. The easy way to find optim tree is to just count the number of rows in a leaf node that are not in the majority class, and d total number of rows. So go through every box, find the smaller number, and add it to every c divide by the total number of rows classified.
- Pessimistic error thinks generalization error is the sum of training error and a penalty term for means a more complex model will result in a higher error. This is done to comply with Occam's

$$-\sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

$$\cdot | t)$$

entropy of just all the

ore prune the tree, or have
er. Typical conditions
e is less than a user
oved from the last node

overfitting. We do this in
arity class, and subtree
A good way to test when a
e things into just one leaf

ake any significantly

a. So there has to be a way
wn rows. There are three
l) known as optimistic

the decision tree created by
mistic error in a decision
divide that amount by the
other smaller number and

r model complexity. This
s Razor, which says that

simpler model is better. To do this, when counting the number of error cases in your leaf nodes total number of rows, add on this amount: (penalty multiplier * number of leaf nodes).

- Validation set error is one way to measure generalization error. It is also called reduced error. separate validation set to measure the error. It just uses a separate set to compute error value validation and 2/3 for training the model. This is the best way to approach generalization error validation set in the optimistic mode. There are a couple of ways to make validation sets:
 - o Single holdout just reserves 2/3 of a set of data for training and 1/3 for testing, or variation
 - o Random subsampling just repeats single holdout k times with different selections of rows
 - o Cross validation partitions data into k disjoint subsets, trains on k-1 of those partitions and It then repeats, ensuring that all k subsets are used as training sets. It uses the averages of these sets to predict the eventual generalization error.
 - o Bootstrapping samples with replacement from the original dataset, with each bootstrap having same size as the original dataset.

ROC Curves

- The precision of a classifier refers to its ability to correctly identify actual positive classes as positive. If you have many false positives, you have low precision.
- The recall of a classifier refers to its ability to correctly classify positive examples as positive, meaning it correctly identifies them. If you have a lot of false negatives, you have low recall.
- The F measure is the harmonic mean of recall and precision. It is a good way to build a strong classifier.
- A ROC curve plots the percentages of your True Positive Rate on the y axis and False Positive Rate on the x axis. You want the ROC curve generated by your classifier to hang towards the upper left corner.
- To draw a ROC curve, the classifier needs to output the likelihood that a specific record will be positive.
- The steps are as follows:
 - o Sort the likelihoods that the classifier outputs from lowest to highest.
 - o Start with the lowest record. Assign it and every record above it to the positive class, and the records below it to the negative class. Compute the true positives, the false positives, the true negatives, and the false negatives. Next, obviously, the true positive rate and the false positive rate as the first point on the curve.
 - o Next, move up to the next column in the table. Classify the new column as positive, as the columns to the right of it as positive, and the columns to the left of it as negative. Recompute TP, FP, TN, FN accordingly. Repeat this until you hit the highest likelihood of positivity record column.
- And that's how you do it.

Logistic Regression

- A logistic regression takes on the form of:

$$\ln \left(\frac{p(X)}{1-p(X)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n$$

es, before dividing by the

pruning, which uses a
es, so like 1/3 of the set for
r. Also, you run the

tions of those proportions.

vs
nd tests on the remaining.
of the error on these test

o validation set having the

positive. If you have a lot of

ot negative. If you have a

classifier.

Rate on the x axis. You

e positive.

d then compute the true
y calculate the true positive

well as everything to the
N, FN and graph the rates

$$\text{logit}(1 - p(X)) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$p(X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

- It is a function that is created based on a set of data. β_0 is a variable that is estimated based on the intercept. β_1 and subsequent variables are also estimated, and these are labeled based on the data.

KNN

- KNN works as follows:
 - o Take a single row of data
 - o Compare its distance using a proximity measure to every other point in your model (which has many rows)
 - o Find the k nearest points, and the class of the new row is the class of the majority of those neighbors.
 - o You can also weight the distances by taking the distance, squaring it, dividing 1 by it, and then summing up all the weights to get the largest value instead of the smallest value when looking at the k neighbors.
 - o KNN are lazy, they do not build models explicitly, they just use training data. They are supervised learning models, but they are based on local info instead of a built global model.
 - o A good value for K is usually around 10 or so. Small K overfit, large k underfit.

Rule Based Classifiers

- A rule based classifiers is a set of rules r that are comprised of implications. The antecedent states one or more attribute tests, and the consequent is always a class.
- We say a rule covers a row if the antecedent matches that row, or the attributes in the row make the test true.
- The coverage of a rule is equal to the fraction of records in the dataframe that trigger the rule AND are correctly classified by the rule, with the denominator being the total number of records.
- The fraction of records that trigger the rule AND are also correctly classified by the rule, with the denominator being the rows the rule covers.
- A good ruleset is mutually exclusive, and this occurs when no two rows trigger the same rule. This means each row triggers at least one rule. Combining these, we ensure that every row is covered by at least one rule.
- If we cannot accomplish this, we use the default rule, which is assigned to the rows that do not trigger any rules in the ruleset.
- If a ruleset is not mutually exclusive, we can deal with this in two ways:
 - o Ordering the rules, meaning the first rule triggered assigns the class.
 - o Using the rule classes as a vote, assigning the row the class with the most amount of rules that triggered it.
- The Sequential Covering method algorithm is as follows:
 - o Start from an empty rule
 - o Grow a rule using the Learn One Rule function

on the data called the
re attributes in the data.

which is just another set of

those k points.

and then looking for the

usceptible to noise, as they

ates a conjunction of

make the antecedent true.

. The accuracy is the
denominator of the fraction

It is also exhaustive,
covered by exactly one rule.
not trigger any of the other

e votes.

- GROW A RULE USING THE LEARN ONE RULE FUNCTION
- Remove training records covered by the rule
- Repeat steps 2 and 3 until you meet stopping criteria
- The LOR function attempts to extract a rule that covers many positive examples and few negative examples
- You can either start with a really specific rule, or a really general rule as your empty rule.
- There are a couple of ways you can evaluate a rule:

Metrics:

- ❖ Accuracy

$$= \frac{n_+}{n}$$

Accuracy does not account for the coverage of the rule!

- Laplace

$$= \frac{n_+ + 1}{n + k}$$

- ❖ M-estimate

$$= \frac{n_+ + kp}{n + k}$$

- Typically the prior is given in most questions.
- Two ways to know how to stop is to set a stopping criterion by measuring gain on accuracy/lambda. If gain is not significant you discard the rule. The other is rule pruning, where you remove one or more rules and compare the error rate on a validation set before and after pruning. If error is better, prune.

Naïve Bayes

- Bayesian Classifiers take a statistical approach to classification. It considers each attribute and its value independently of the other variables. Given a record with attributes (x_1, x_2, \dots, x_N) , the goal is to predict class C. Specifically, we want to find the value of C that maximizes $P(C | X_1, X_2, \dots, X_N)$.
- The approach is to use Bayes Theorem to predict the probability of the class given all the attributes. To do this, we need to find:
 - The probability of all the attributes together given a class $P(X | C)$
 - Times the probability of the class $P(C)$

tive examples.

r

stances

positive

l by rule

sses

ty

place/m-estimate. If the
f the conjuncts in a rule,
ne.

l class label as random
ally, we want to find the

bute values, and to find

- With all of this divided by the probability of all the attributes. $P(X)$
 - $P(X)$ gives kind of a scaling factor or normalizing factor, not very interesting. We are interested in the relative probabilities.
- The primary question is how to estimate $P(X | C)$. Bayes assumes the variables are independent given C , so that $P(X|C)$ can be factored into a product of multiple, per attribute likelihoods. EX: $P(X_1, X_2, \dots, X_n | C) = P(X_1 | C_j) P(X_2 | C_j) \dots P(X_n | C_j)$.

$$P(X_1, X_2, \dots, X_n | C) = P(X_1 | C_j) P(X_2 | C_j) \dots P(X_n | C_j)$$

❖ Can separately estimate $P(X_i | C_j)$ for all X_i

- This algorithm is naïve since most attributes in an attribute vector are not independent of each other. There is some correlation between attributes in an attribute vector. Bayes throws that out the window.
- For discrete attributes, $P(X_i | C_i) = \frac{\text{number of instances of } X_i \text{ in class } C_i}{\text{total number of instances in dataset}}$ = number of instances of rows having attribute X_i and class C_i divided by total number of instances in the dataset times the class C_i appears in the dataset.
- If you have a continuous attribute, you can discretize the range into bins, like low, medium and high, or estimate a probability distribution density estimation, by assuming the attributes are normally distributed (or some other distribution), and use the data to estimate the parameters for the normal distribution for that attribute. Then we can use the formula for the normal distribution to estimate the probabilities of our continuous attributes. This is done using the sample variance of the data, which is plugged into this equation that outputs a probability:

$$P(X_i | C) = \frac{1}{\sqrt{2\pi}\sigma} \exp^{-\frac{1}{2} \frac{(X_i - \mu)^2}{\sigma^2}}$$

- For example, if you have a row with three attributes, then $P(X|C) = P(\text{first attribute} | C) * P(\text{second attribute} | C) * P(\text{third attribute} | C)$.
- The slide 125 skips the step denoting the fact that $P(X|\text{Class} = \text{No})$ is equivalent to $P(\text{Class} = \text{No} | X)$.
- We run into an issue when one of the probabilities in the algorithm is zero. To fix this, we do Laplace smoothing. We add one to the number of attributes of type X_i that belong to class C_i , and adds the number of classes in the data set to the denominator, which if you recall is the number of classes in the dataset N . Here is the formula:

$$\text{Original: } P(X_i | C) = \frac{N_{ic}}{N}$$

interested in $P(X|C) * P(C)$.

nt given class. This means

$$\dots, X_N | C) = P(X_1 | C) *$$

$$P(X_n | C_j)$$

and C_j .

h other. There is generally

v.

i divided by the number of

d high. You can also

y distributed (or some

or each class, and then use

the sample mean and the

cond attribute | C) *

$\dots | X_1, X_2, \dots, X_N)$.

Laplace Smoothing, which
of records with the class C_i
total. The equations are

c: num

p: pri

m: pa

- Laplace: $P(X_i | C) = \frac{N_{ic}+1}{N_c+c}$

m-estimate: $P(X_i | C) = \frac{N_{ic}+mp}{N_c+m}$

- Bayesian Classifiers have a goal to choose a value of C such that $P(X_1, X_2, \dots, X_N | C) * P(C)$ is maximized. There exists to justify this maximization, that basically says that the maximum probability of class creates the best classifier.
- Bayes decision boundaries on simulated training data describe the lowest possible test error rate. This is because the maximum classifier probability/accuracy of classifier. Your model running on the test data can have a lower than the one that the BDB describes.
- In the real world, we do not have the ability to create the BDB as we cannot compute the actual probabilities of all features of a collected dataset. This is why we need to come up with estimates for the true and unknown probabilities. We hope that the algorithm that we use to estimate $P(+ | x)$ from the data is close to the true probability, so we can have a pretty solid test error rate.
- Don't worry too much about Bayesian Classifiers, he skipped the example slide.

Artificial Neural Networks

- An ANN is an interconnected assembly of nodes and directed edges that can be used to learn complex functions between a set of inputs and outputs. The arrangements of these nodes and edges is called the architecture of the network.
- A ANN is made up of perceptrons linked together. The input node of the perceptron transmits its input through the outgoing link without any transformation. The output node is a mathematical device that takes the sum of its inputs, subtracts the bias term, and then produces an output that depends on the activation function it chooses. Here is the equation for the activation function of a sign AF:

number of classes
or probability
parameter

maximized. A theorem
states the most accurate

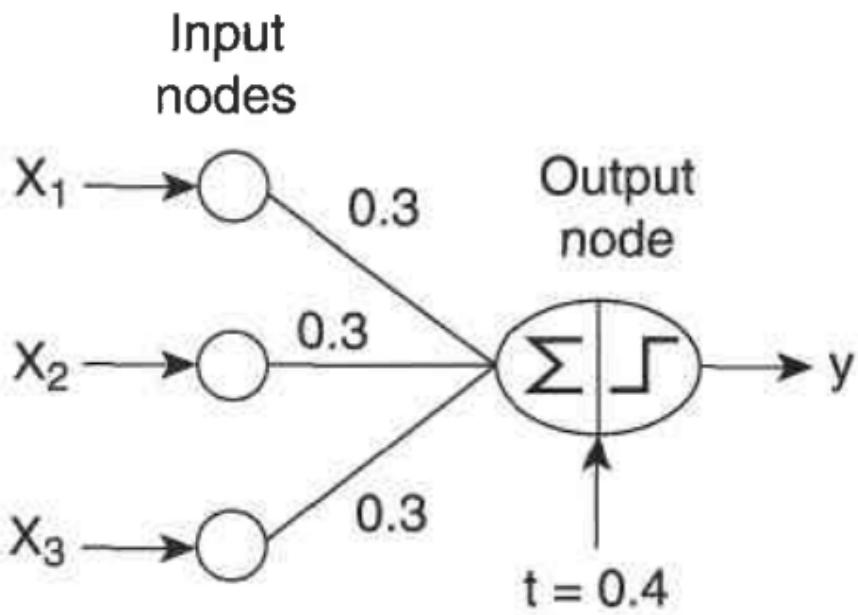
rate, which is based on the
not have an error rate

ual posterior probabilities
own probabilities from the
e true posterior

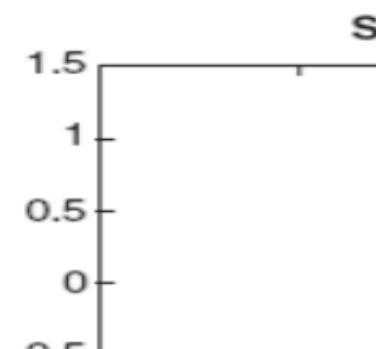
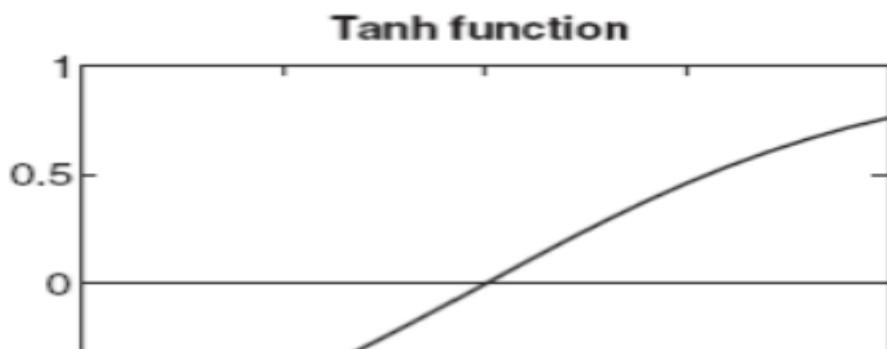
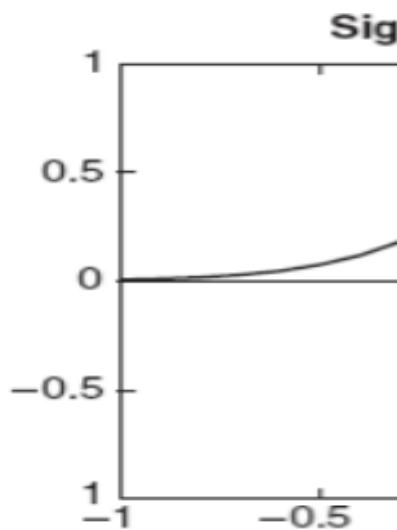
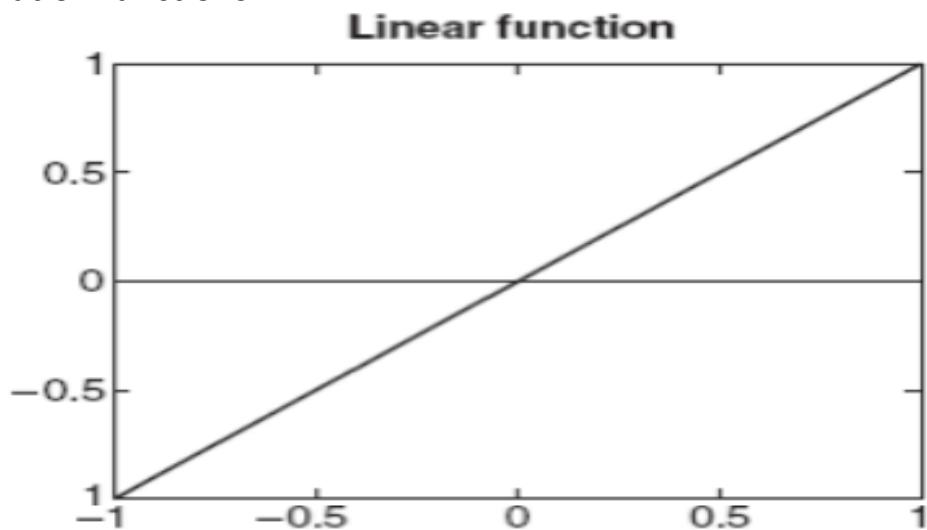
the functional relationship
e architecture of the

s the value it receives to
computes the weighted
ign of the resulting sum. It
uation for the activation

- $y = \text{sign}(w_d x_d + w_{d-1} x_{d-1} + \dots + w_2 x_2 + w_1 x_1 + b)$
- Where w is the value of the weights that deliver the x attribute to the perceptron, and t is the threshold.



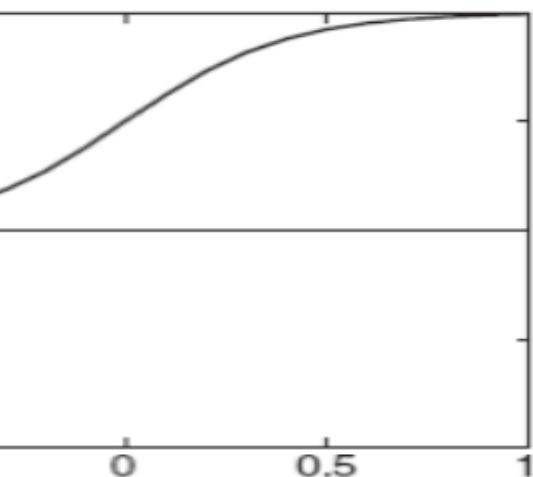
- Activation functions:



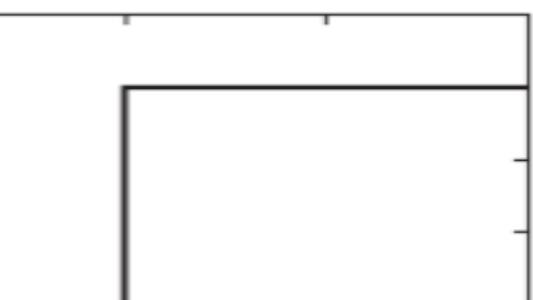
$$v_1x_1 - b)$$

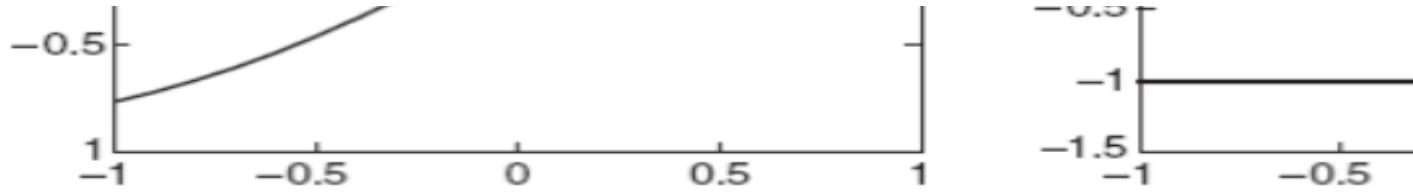
bias term:

Sigmoid function



Sign function





- The learning algorithm for an ANN is as follows:
 - o Give each input weight a random value
 - o For every training example, compute the predicted output.
 - o For every weight, update the weight using the weight update formula:

$$w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_i$$

- Where the new weight, w_{jk+1} , is the old weight plus the user specified learning rate times the times the actual class of the row minus the predicted output of the perceptron.
- A single perceptron only separates linearly separable data. It will not converge with data that is not linearly separable.
- To get around this, an ANN is constructed with many perceptrons, with an initial layer being the input layer, a hidden layer being a bunch of perceptrons in the middle, and the output layer outputting the final y value. In this network, you adjust the weights in a way that minimizes the objective function so the ANN can correctly predict the class labels of your training examples:

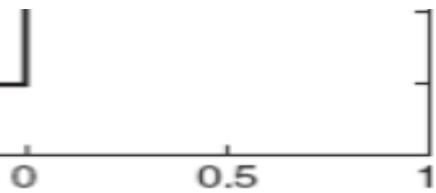
Objective function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

- Backpropagation has two phases. In forwards phases, the weights obtained from the previous phase are used to compute the output value of each neuron in the network, starting at the input level and working through the layers one at a time. During the backward phase, the weight update formula is applied in the reverse direction, starting from the output node and moving backwards to the input layer. All of this is done to find the weights that minimize the objective function.

Support Vector Machines

- SVM is a function that finds a dividing line between data in a dataset. The line that maximizes the margin between the two closest data points of two classes with a dividing line in the middle is the best line and the decision boundary.



j ,

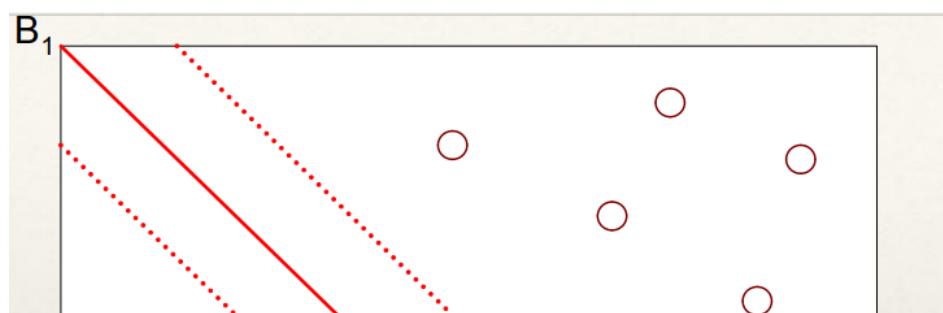
the input for that weight

is not linearly separable.

the input layer, the hidden
value it predicts. To train
N is consistent with the

s iteration are used to
ing to the output, one level
ection, starting at the
hat minimize the output of

the margin between the
one the model attempts



to find.

- B1 is better than B2.
- Rows less than the margin are squares, rows greater than the margin are circles, and if anything between b11 and b12, you can use slack variables to tolerate small errors in the training data.
- If your data's decision boundary is not linear, then you will have to transform the data into a higher dimension where the data is linearly separable.
- The kernel trick is what can do this without the curse of dimensionality. This trick basically represents attributes via their N x N pairwise similarity values. The function that computes the similarity, so that it represents a dot product in some high dimensional feature space, but is easily computed.

Ensemble Classifiers

- Ensemble Classifiers construct a set of classifiers from the training data, and predict class labels for new records by aggregating the predictions made by multiple classifiers.
- If we have 25 base classifiers that run on a row, and we pick the majority vote class, AKA the class predicted by the most classifiers. Another way is consensus, where the classifiers all have to be the same or the error rate is zero. We assume independence of the classifiers.
- If we have an actual positive case, and 13 or greater of the classifiers predict a negative, this is a false negative.
- The probability that ensemble classifiers make a wrong prediction:

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06$$

- The error rate is the epsilon, and this uses an epsilon of 0.35. If the epsilon is greater than 0.5, the classifier will perform worse. Also, if you use the same classifier 25 times, it won't do anything better than 0.35.
- Ensembles work as different classifiers can make up for the negative aspects of other classifiers. The correlation between base classifiers improves the performance of this method.
- Bias-variance decomposition is a formal method for analyzing the prediction error of a predictor. It is the sum of three components: Variance squared + noise.
 - o Bias is the difference between the average prediction and the actual value. Models with high bias do not capture the regularities in the data.
 - o Variance measures the spread of the predictions around the average prediction given similar training data. High variance methods overfit training data.
 - o Bias and variance are inversely related to each other in general. We want low bias and low variance.

ng exists in the margin
higher dimensional space
presents the input
called the kernel, is chosen
uted.

ls of previously unseen

class that is returned by the
e class remains unknown.

s an error.

, then the ensemble
g, the error rate will remain

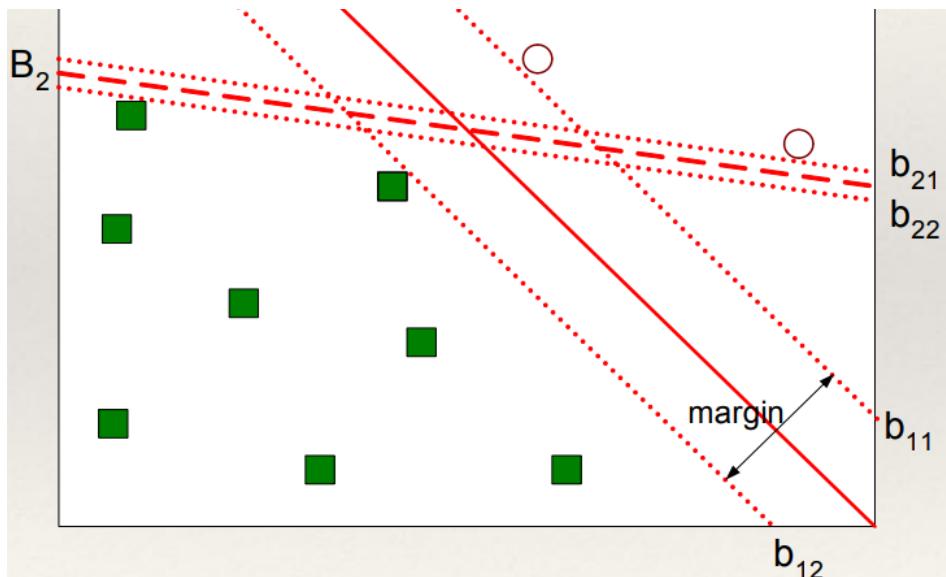
rs in the ensemble. Low

tive model. Error = Bias +

n a high bias underfit, and

small perturbations in the

low variance, with error



being just our noise, things we cannot control.

- Jason uses the term base classifier, which are the classifier algorithms used in the ensemble classifier. We sample from the original training data and use an unstable learner, with the most common being decision trees. We want to pick classifiers with high variance, as combining these will reduce overall variance. These classifiers have low bias.
- Unstable just means we have a high variance learner. This means if we change the training data slightly, the model will be changed dramatically.
- We can construct an ensemble of classifiers in a couple of different ways:
 - o Manipulating the training dataset (bagging)
 - o Manipulating the input features (random forest)
 - o Manipulating the learning algorithms for the classifiers
 - o Manipulating the class labels (Usually done on large class number rows)
- In bagging, you randomly select row attributes with replacement multiple times and build classifiers. We follow the standard ensemble instructions.
- On the 195 training set, the best binary decision stump can perform at a 30% error rate.
- Boosting is an iterative procedure to adaptively change distribution of training data by focusing on misclassified records.
 - o Initially, all N records are assigned equal weights, but these may change at the end of a boosting iteration.
 - o Records with wrong classifieds increase weights, right classifieds decrease weights
 - o Weights are used either to change the sampling distribution in each iteration or to bias the classifier.
- Adaboost takes a base set of classifiers, and weights the importance of a classifier according to how well it actually is.
- The simple algorithms typically have low variance, with high bias. These are logistic regression, linear regression, naïve Bayes, and rule based. Complex typically are more flexible, high variance, and low bias which can lead to overfit.

Review Lecture

- Appendix material can be on the exam.
- Remember data mining is the non trivial extraction of implicit, previously unknown and potentially useful information from data. It is a blend of data analysis methods with sophisticated algorithms for processing large amounts of data.
- We divide tasks into prediction and description. Before the midterm, we talked about prediction methods like regression, classification, and deviation detection. After, we'll cover descriptive methods like clustering and association rule mining.
- Proximity measures are really important!
- His highlights from chapters 2 and 3:
 - o Attribute type: nominal, ordinal, interval, ratio
 - o Know multidimensional density and distributions with the covariance matrix, and the correlation coefficient between attributes. Have a good understanding of how to understand what the correlation is in a dataset.
 - o Know Bayes Theorem. Maybe just write it down.

classifier. Most ensembles
being decision trees. So we
unstable learners should

ta even just a little bit, the

ssifiers on them. Then you

ng more on previously

boosting round.

the underlying model
o how high its error rate

n, KNN, decision trees,
which makes them tend to

tially useful information
large volumes of data.
ion methods, classification,
association.

orrelation value between
a graph.

- Proximity measures, again. Know the 1D views, and the measures for multiple dimensions and Minkowski. SMC, Jaccard, Cosine. Jaccard is better for sparse datasets, as it does not can skew your proximity measure. Cosine is useful for nonbinary data, where attributes
- One thing that could be on the exam is to get a box or scatterplot and interpret something comfortable with all the visualization tools.
- Data quality: noise/outliers, missing values
- Principle Component Analysis is fair game. Interpretation of data, or do part of the calculation you do PCA, what does it mean, what it is it trying to accomplish.
- Dimensional Reduction creates a new, reduced set of features from the original. Like going from more descriptive attributes. Feature subset selection just takes a chunk of the attributes with little to no transformations performed on them.
- His highlights from chapters 4 and 5:
 - The purpose of classification is to learn a model from your training data that generalizes to new data. Overfitting occurs when your model works well on your training set but only on your test set. Underfitting occurs when your model doesn't even fit your training set correctly and give any worthwhile description of your data.
 - You want a good balance point between the overfit and underfit of your model. You want to minimize error.
 - Producing a class label can be broken down into two steps. For a given record, you output a probability or "score" for each possible class label. So a model takes in a record, and outputs these probabilities to predict the class with a varying threshold t between 0 and 1. If $P(+|X) \geq t$, then $x = +$. If $P(-|X) \geq t$, then $x = -$. If $P(+|X) = P(-|X) = 0.5$ which is the default decision threshold, you output a confusion matrix. By altering the threshold, you can build a ROC curve, and output a confusion matrix for each value of t .
 - By increasing the threshold, you typically decrease TP and FP, and increase FN and TN. Underfitting occurs when the threshold is too low. If the threshold is too high, precision goes up when increasing t , you have a good model. This means that your model is correctly classifying most of the positive assignments, as most TPs that are strongly believed to be positive will lie above high thresholds.
 - Have a good idea of how decision trees, naïve Bayes, rule based and KNN work. Do not worry about neural networks, SVM, and ensembles.

Exam Answer Notes

- An equal interval width is when discretizing a continuous variable where you create each interval at the same width, like 10, 20, 30, 40. An equal frequency distribution is when the same amount of data points are binned into each interval.

ns. He mentions Euclidian
ot count 00 matches which
are not T/F or 0/1.
ing from that. Be

ulation. Understand how

ing from 20 attributes to 4
s and only uses them with

s well to previously unseen
our training set, it doesn't
ectly, it is too general to

nt to minimize your test

ut a posterior probability
cores. The second step is
. Otherwise, $x = -$. When t is
 t incrementally, you can

Upon increasing t , if your
el is confident in its class
values.
worry about deep details

erval of discrete type at an
ntinuous variables get

