

Classification is the process of categorizing data into predefined classes or categories. For

## Chap-8

### ① Classification (বিভাগ কৰা)

### ② " " (কী?)

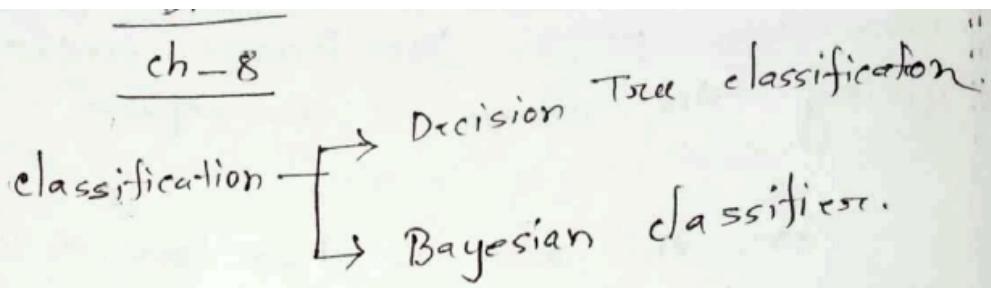
Classification is the process of categorizing data into predefined classes or categories.

For example, email spam filters use classification to automatically label incoming emails as either "spam" or "non-spam" based on learned patterns.

### ② " " (কী?)

Classification is sorting data into predefined categories. For example, spam filters categorize emails as "spam" or "non-spam" based on learned patterns.

ch-8



classification → কোন data কোন class এ পাই।

মেলা: যুক্তিগত dataset ন হলে loan নিল  
safe তা risky zone ন আছি এটা  
classify করা হবে।

8.1.2 → General Approach

data দিয়ে machine কে learn করার পথে—  
যদি মনে নেওয়া data আসলে তবে সেটা পার্স  
সে আসলে হচ্ছে কী!

data tuple কে যান, sample, example,  
data point, object .

- classification ২ ধরণ কোন সাথে →  
1) supervised Learning  
2) unsupervised "

④ Tree pruning কি হচ্ছে কী?

\* pruning  $\rightarrow$  unnecessary branch remove  
 $\rightarrow$  accuracy  $\rightarrow$  what?

→ Prepruning ~~Approach~~ for 3 methods  
→ Post pruning ~~for~~  $\rightarrow$  for?

Prepruning is a strategy in decision tree algorithms where the growth of the tree is restricted during training to prevent overfitting. It involves setting conditions, such as maximum depth or minimum samples, to control the complexity of the tree from the start.

→ Post pruning ~~for~~  $\rightarrow$  for?

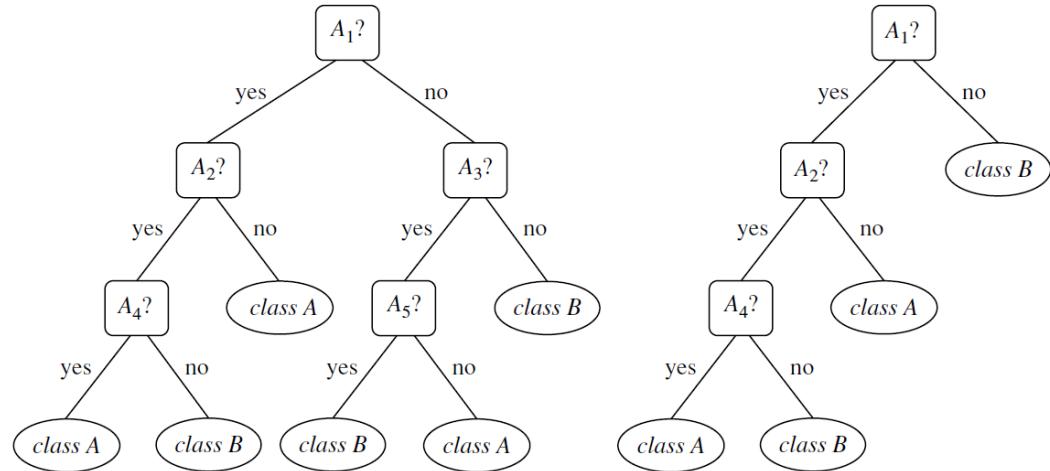
Postpruning is a technique for decision trees where, after the tree is fully grown, certain branches are removed or collapsed to enhance generalization on new data. This helps prevent overfitting and improves the model's performance.

*“How does tree pruning work?”* There are two common approaches to tree pruning: *prepruning* and *postpruning*.

In the **prepruning** approach, a tree is “pruned” by halting its construction early (e.g., by deciding not to further split or partition the subset of training tuples at a given node). Upon halting, the node becomes a leaf. The leaf may hold the most frequent class among the subset tuples or the probability distribution of those tuples.

When constructing a tree, measures such as statistical significance, information gain, Gini index, and so on, can be used to assess the goodness of a split. If partitioning the tuples at a node would result in a split that falls below a prespecified threshold, then further partitioning of the given subset is halted. There are difficulties, however, in choosing an appropriate threshold. High thresholds could result in oversimplified trees, whereas low thresholds could result in very little simplification.

The second and more common approach is **postpruning**, which removes subtrees from a “fully grown” tree. A subtree at a given node is pruned by removing its branches and replacing it with a leaf. The leaf is labeled with the most frequent class among the subtree being replaced. For example, notice the subtree at node “ $A_3?$ ” in the unpruned



**Figure 8.6** An unpruned decision tree and a pruned version of it.

⑤ scalability - BOAT Approach (347)

→ DT for real-life data set 1m, 2m रियल डेटा सेट, जो बड़ा है

(2)

मेमोरी स्पैटिक मल्लांगे. SO BOAT  
Approach statistical measure ~~of~~ sample रिहाई

विभिन्न ट्री बनाएं एवं ट्री के अनुभवों  
एवं विभिन्न सैमप्ल बनाएं।  
Subset इन ट्री के बारे में बताएं।  
जटिलता - voting Technique बताएं।  
ट्री गणना जोड़कर एवं बनाएं।

→ multiple root, trigger rule वा APPROACH 9  
size order, rule ordering class base (358)

The size ordering scheme assigns the highest priority to the triggering rule that has the "toughest" requirements, where toughness is measured by the rule antecedent size. That is, the triggering rule with the most attribute tests is fired.

→ COVERAGE, ACCURACY (356)

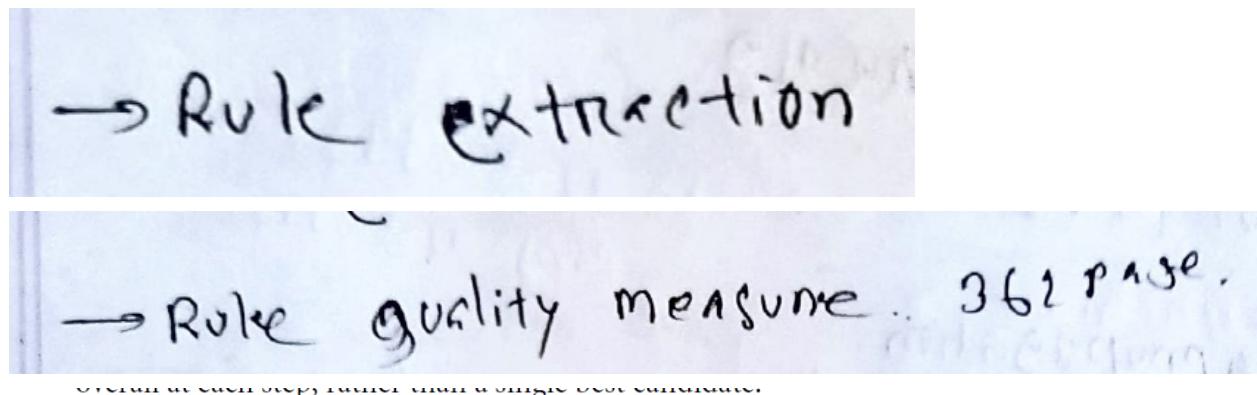
If the condition (i.e., all the attribute tests) in a rule antecedent holds true for a given tuple, we say that the rule antecedent is **satisfied** (or simply, that the rule is satisfied) and that the rule **covers** the tuple.

A rule  $R$  can be assessed by its coverage and accuracy. Given a tuple,  $X$ , from a class-labeled data set,  $D$ , let  $n_{covers}$  be the number of tuples covered by  $R$ ;  $n_{correct}$  be the number of tuples correctly classified by  $R$ ; and  $|D|$  be the number of tuples in  $D$ . We can define the **coverage** and **accuracy** of  $R$  as

$$\text{coverage}(R) = \frac{n_{covers}}{|D|} \quad (8.16)$$

$$\text{accuracy}(R) = \frac{n_{correct}}{n_{covers}}. \quad (8.17)$$

That is, a rule's coverage is the percentage of tuples that are covered by the rule (i.e., their attribute values hold true for the rule's antecedent). For a rule's accuracy, we look at the tuples that it covers and see what percentage of them the rule can correctly classify.

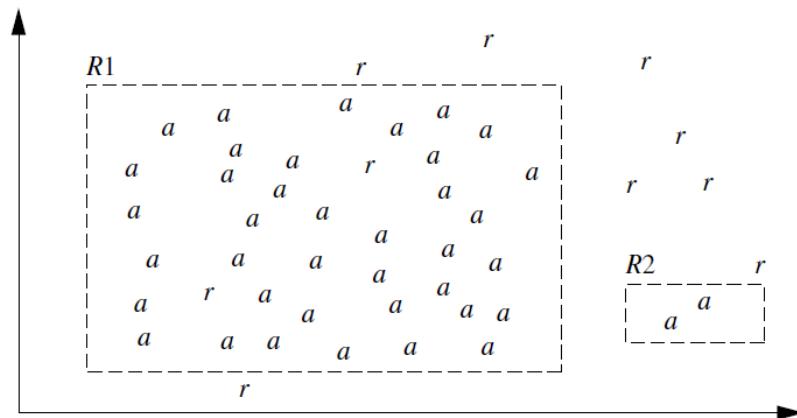


## Rule Quality Measures

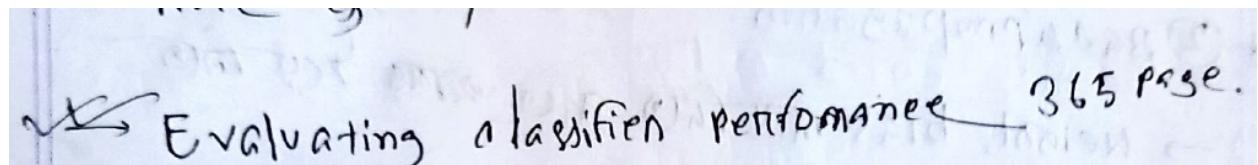
*Learn\_One\_Rule* needs a measure of rule quality. Every time it considers an attribute test, it must check to see if appending such a test to the current rule's condition will result in an improved rule. Accuracy may seem like an obvious choice at first, but consider Example 8.8.

**Example 8.8 Choosing between two rules based on accuracy.** Consider the two rules as illustrated in Figure 8.12. Both are for the class  $loan\_decision = accept$ . We use “ $a$ ” to represent the tuples of class “*accept*” and “ $r$ ” for the tuples of class “*reject*.” Rule  $R1$  correctly classifies 38 of the 40 tuples it covers. Rule  $R2$  covers only two tuples, which it correctly classifies. Their respective accuracies are 95% and 100%. Thus,  $R2$  has greater accuracy than  $R1$ , but it is not the better rule because of its small coverage. ■

From this example, we see that accuracy on its own is not a reliable estimate of rule quality. Coverage on its own is not useful either—for a given class we could have a rule that covers many tuples, most of which belong to other classes! Thus, we seek other measures for evaluating rule quality, which may integrate aspects of accuracy and coverage. Here we will look at a few, namely *entropy*, another based on *information gain*, and a *statistical test* that considers coverage. For our discussion, suppose we are learning rules



**Figure 8.12** Rules for the class  $loan\_decision = accept$ , showing *accept* ( $a$ ) and *reject* ( $r$ ) tuples.



8.5.21

### Model evaluation

Accuracy  $\rightarrow$  recognition scale  $\rightarrow$   $\frac{1}{2}$   
Model evaluation  $\rightarrow$  under  $\rightarrow$ .

posi page 364

positive plus / Yes class / positive tuple  
positive negative " / No " / Negative tuple.

Multiple class  $\rightarrow$   $\frac{1}{2}$   $\frac{1}{2}$   $\frac{1}{2}$  class  $\rightarrow$   
target target yes  $\rightarrow$  class  $\rightarrow$   
actual No class.

8.14

P  $\rightarrow$  actual dataset  $\rightarrow$  positive

P'  $\rightarrow$  " Model  $\rightarrow$  positive  $\rightarrow$

N  $\rightarrow$  " actual dataset  $\rightarrow$  negative

N'  $\rightarrow$  " Model  $\rightarrow$   $\rightarrow$

$$\frac{6054 + 2588}{10000} = 0.8542$$

$$\therefore \text{accuracy} = \frac{TP + TN}{P + N}$$

$$\text{Error Rate} = \frac{FP + FN}{P + N}$$

$$1 - \text{accuracy}$$

class imbalance problem → याचा data ती imbalanced  
वा biased हस्ते आका data ती.

मानव: 100 वरील class असू असू येचा class आव  
3 देऊ No class. तिने अगाडी असू असू 97%. accuracy  
हस्ते, तिने model मध्ये 3 देऊ No class येत असू  
आव, तिने accuracy असू आव ता। आशल  
आव, आशल accuracy असू आव ता। आशल  
आव मासू class imbalance problem नाही accuracy  
असू matrix ता। असू आव ता -

$$\text{sensitivity} = \frac{TP}{P} \rightarrow \text{आव true positive ratio असू आव  
indicate असू।}$$

आव positive tuple असू आव कोरेक्टly  
identity आव आव कोरेक्ट

$$\text{आव specificity} = \frac{TN}{N} \rightarrow \text{आव negative}$$

tuple द्वारा कैसे सही रूप से आवाहन

$$\text{accuracy} = \text{sensitivity} \frac{P}{P+N} + \text{specificity} \frac{N}{P+N}$$

→ confusion matrix

→ sensitivity, specificity. ③ 6 7

page 368

sensitivity অথবা specificity প্রায় same  
 একই আভ্যন্তর বলতে পাই সে model এর অভ্যন্তর  
 আলো performance করতেছে। সবচেয়ে মুক্তি আন্তর  
 আলো আসে তখন আভ্যন্তরীণ প্রযোগ আবশ্যিক  
 এটি accurate scenario না।

Precision  $\rightarrow$  measure of exactness  
 recall  $\rightarrow$  " completeness

$$\text{Precision} = \frac{TP}{TP+FP} \rightarrow \text{আভ্যন্তরীণ model}$$

$\downarrow$

কৃতিত্ব করতে পারে

A value  $1-23\% / 100\%$  positive বলে।

আলো আভ্যন্তরীণ class  $\rightarrow$  yes class

মনে রাখো  $\text{belong } \text{positive}$

প্রযোগে  $\text{belong } \text{positive}$  করতেছে।

recall  $\equiv 1 / 100\%$ . यहां आज FN का तो  
काला term नाही तो आज जो आवृत्ति असली  
येस  $\rightarrow$  model असली हो ना रहती ही।  
but इस extra किस data को इस तरह  
माझे positive घटी याऊन ठिक हो।

आखे precision वा recall आजकल म्हणा  $F_1$   
होतो ही। आखे याची combined matrix  
आशयांची दोनों ओर  $F_1$  score.

$$F_{\text{mea}} = \frac{\text{Harmonic Mean} = \frac{1}{\alpha} + \frac{1}{\beta}}{\text{Weight factor precision} \downarrow \quad \downarrow \text{recall}}$$

$$F = \frac{(2 \times P \times R)}{P + R}$$

Scalability  $\rightarrow$  large dataset याचे काढ यावा  
परफॉर्म होता ही।

3\* → Rock curve - 375 page. (~~AUC, TP, FP, TPR, FPR~~)  
sort नसीहत,  $\dots$

ROC curve  $\rightarrow$  410 page

cost benefit analysis  $\Rightarrow$  1

$\rightarrow$  cost associated with false negative is greater than false positive. ex: cancer patient 100 रुपयों में 3 रुपये खर्च करता है 95 रुपये non-cancer.

ROC curve  $\rightarrow$  multiple threshold  $\Rightarrow$  check  $\Rightarrow$  threshold

मात्र:

	yes	no
yes	TP	FN
no	FP	TN

$\Rightarrow$  TPR

$$\text{sensitivity} = \frac{TP}{P}$$

$$\text{specificity} = \frac{TN}{N}$$

$\Rightarrow$  overall accuracy  $\Rightarrow$   $\Rightarrow$  मात्रा मात्रा  
मात्रा मात्रा mathematically  
 $\Rightarrow$  threshold

~~1 - TN/N~~

$$1 - \frac{TN}{N} = \frac{FP}{N} = FPR$$

Page 412 \*

Prob  $\rightarrow$  0.90 आने 0.9 यह तो इसके positive

अपने 0.9 " यह n Negative.

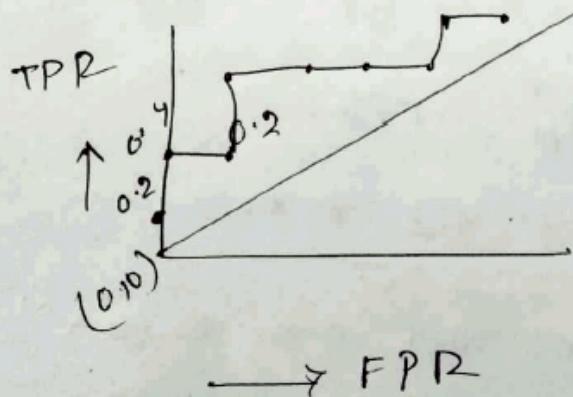
\* मर्गीदार सूची का नित अप्रेस probability (G)

			TP	FP	TPR	
1	P	0.9 (P)	1	0	1/1 = 1.0	1.0
2	P	0.8 (P)	2	0	2/5 = 0.5	0
3	N	0.7 (P)	2	1	0.7	0
4	P	0.6 (P)	3	1	0.4	0.2
5	P	0.55 (P)	4	1	0.6	0.2
6	N	0.54 (P)	4	2	0.8	0.2
7	N	0.53 (P)	4	3		
8	N	0.51 (P)	4	4		
9	P	0.50 (P)	5	4		
10	N	0.40 (P)	5	5		

wrong

ROC curve

(11)



curve मानवी  
0-1 आणि एवा  
आपले probability  
नियंत्रण दिला

Ex. 8. 12

TP 2 मध्यात 1 TN 3 मध्यात 1 ...

## ⑧ Ensemble methods.

→ Bagging, Boosting, Ada Boost, - frenet, 0.1 class imbalance

→ Resample data? ①

→ Random forest.

→ oversampling ③ ④

→ under sampling ③ ④ page 8.12 example

→ Threshold moving.

## Ensemble Learning

(i) Bagging

(ii) Boosting

(iii) Random Forest

## Boosting

Adaboost মাত্রে adaptive (

DM  
8.6 → 414 Page

08/12/23

classification improve করে ensemble model  
ধূম্র ।

majority voting টির ম্যানে  
calco ১৪১ → ২৪ । ৫০% হল  
average accuracy ৭৫% ।

## Bagging

আন্তর্বর্তীর্থে voting technique. Parallel  
process হয়ে ।

Boosting → sequentially হয়ে

গোপনীয় model কর ও weight সম্পর্কীয় model র উপর  
বর্তো ।

$t_1 \rightarrow w_1$   
 $t_2 \rightarrow (w_2) \rightarrow w_1$  এর পৃষ্ঠা depend  
— হাতে ।

→ ~~मात्र~~ मात्र hi model first train करे यही first iteration  
एवं करते ही उसे उपर्युक्त weight दिया जाता है first iteration में  
यह अन्य उदाहरण का उपर्युक्त weight update नहीं  
करता आगे बढ़ते ही modify करता है।

→ slide पर इसी प्रक्रिया चलते हैं।

### AdaBoost

→ First एवं tuple के weight same होंगे  
like:  $\frac{1}{d}$

→ अगर यह tuple incorrectly classify करता है  
उपर्युक्त weight बढ़ाया जाता है। weight बढ़ाया  
जाता है कि model के classify करना easy.  
weight उपर्युक्त reflect करता है आगे model  
classify करना भी hard.

8.6.4 Random Forest  $\rightarrow$  419 Page

→  
अनेक उपर्युक्त decision tree रख सकते हैं,  
majority voting  $\rightarrow$  ऐसे class आनंदित  
होते हैं।

8·6·5

class Imbalance problem  $\rightarrow$  point of class  
interest class  $\Rightarrow$  into positive class.

→ Traditional algo  $\nexists$  imbalance data  
→ निम्न concerned w/

Oversampling: - 100 " to positive class  
↓  
1000 " Negative "

↓  
ଆଜେ ଅନ୍ତର୍ଗତ point of interest class ଏବଂ  
duplicate କୌଣସି ହୋଇବାକୁ ବ୍ୟାଖ୍ୟାନ କରିବାକୁ 100 ମୀଟର୍ ଦିଶା  
replica କୌଣସି Negative class ଏବଂ ସମ୍ପାଦନ କରାଯାଇବାକୁ ।

अब undersampling हला Negative class को (अफ्पे) eliminate गर्नु चाहिए।

ଭାବରେ dataset ରେ ଆମାର କଲ୍ୟ କୋଣପେ କିମ୍ବା ଆବଶ୍ୟକ କିମ୍ବା ନିତ୍ୟ କବାଳେ oversampling ଏବଂ କିମ୍ବା undersampling କାହାରେ ଦିଆଯାଇଛି ?

The **threshold-moving** approach to the class imbalance problem does not involve any sampling. It applies to classifiers that, given an input tuple, return a continuous output value (just like in Section 8.5.6, where we discussed how to construct ROC curves). That is, for an input tuple,  $X$ , such a classifier returns as output a mapping,

## 8.7 Summary

- **Classification** is a form of data analysis that extracts models describing data classes. A classifier, or classification model, predicts categorical labels (classes). **Numeric prediction** models continuous-valued functions. Classification and numeric prediction are the two major types of prediction problems.
- **Decision tree induction** is a top-down recursive tree induction algorithm, which uses an attribute selection measure to select the attribute tested for each nonleaf node in the tree. **ID3**, **C4.5**, and **CART** are examples of such algorithms using different attribute selection measures. **Tree pruning** algorithms attempt to improve accuracy by removing tree branches reflecting noise in the data. Early decision tree algorithms typically assume that the data are memory resident. Several scalable algorithms, such as **RainForest**, have been proposed for scalable tree induction.
- **Naïve Bayesian classification** is based on Bayes' theorem of posterior probability. It assumes class-conditional independence—that the effect of an attribute value on a given class is independent of the values of the other attributes.

- A **rule-based classifier** uses a set of IF-THEN rules for classification. Rules can be extracted from a decision tree. Rules may also be generated directly from training data using sequential covering algorithms.
- A **confusion matrix** can be used to evaluate a classifier's quality. For a two-class problem, it shows the *true positives*, *true negatives*, *false positives*, and *false negatives*. Measures that assess a classifier's predictive ability include **accuracy**, **sensitivity** (also known as **recall**), **specificity**, **precision**,  $F$ , and  $F_\beta$ . Reliance on the accuracy measure can be deceiving when the main class of interest is in the minority.
- Construction and evaluation of a classifier require partitioning labeled data into a training set and a test set. **Holdout**, **random sampling**, **cross-validation**, and **bootstrapping** are typical methods used for such partitioning.
- Significance tests and ROC curves are useful tools for model selection. **Significance tests** can be used to assess whether the difference in accuracy between two classifiers is due to chance. **ROC curves** plot the true positive rate (or sensitivity) versus the false positive rate (or  $1 - specificity$ ) of one or more classifiers.
- **Ensemble methods** can be used to increase overall accuracy by learning and combining a series of individual (base) classifier models. **Bagging**, **boosting**, and **random forests** are popular ensemble methods.
- The **class imbalance problem** occurs when the main class of interest is represented by only a few tuples. Strategies to address this problem include **oversampling**, **undersampling**, **threshold moving**, and **ensemble techniques**.

- Datas
- ③ SVM - support vector machine 408
- sum (9)  $\rightarrow$  sum Kernel  $\rightarrow$  70
- Theory → linearly separable for 3? → n - separable?

- ④ Lazy Learners vs Eager learners 423
- Naive Bayes for LL or EI?
- KNN for LL?
- KNN 70 math
-

Feature	Eager Learner	Lazy Learner
When model is built	During training phase	During prediction phase
Prediction speed	Faster	Slower
Training speed	Slower	Faster
Memory usage	Lower	Higher
Example algorithms	Linear regression, decision trees	k-nearest neighbors, support vector machines

## Lazy vs. Eager Learning

- Lazy vs. eager learning
  - Lazy learning (e.g., instance-based learning): Simply stores training data (or only minor processing) and waits until it is given a test tuple
  - Eager learning (eg. Decision trees, SVM, NN): Given a set of training set, constructs a classification model before receiving new (e.g., test) data to classify
- Lazy: less time in training but more time in predicting
- Accuracy
  - Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form its implicit global approximation to the target function
  - Eager: must commit to a single hypothesis that covers the entire instance space





Naive Bayes is considered an **eager learner** in machine learning. Here's why:

- **Model building:** During the training phase, Naive Bayes constructs a complete probabilistic model of the data. This model represents the relationships between features and the target variable. Once built, the model is used for prediction without further modifications.
- **Prediction speed:** Once trained, Naive Bayes can make predictions quite quickly on new data points. It simply calculates the posterior probability for each class based on the features of the new data point and chooses the class with the highest probability.
- **Training speed:** While Naive Bayes is relatively fast at predicting, training it can be slower compared to some lazy learners, especially for large datasets. This is because it needs to build the entire probabilistic model upfront.

Therefore, based on the characteristics of eager learners, Naive Bayes fits the bill. It builds a complete model during training and uses it readily for predictions without further computation on new data points.

I hope this explanation clarifies the categorization of Naive Bayes! Feel free to ask if you have any further questions.

**k-Nearest Neighbors (k-NN) is a lazy learner because it doesn't build a model during training. Instead, it memorizes the entire dataset and makes predictions for new instances based on the closest neighbors at prediction time.**

1. **Delayed Model Building:** Unlike eager learners like Naive Bayes that build a complete model during training, kNN doesn't create a global model beforehand. Instead, it stores all the training data points. When given a new data point for prediction, kNN searches for its  $k$  nearest neighbors within the stored data. This "on-demand" approach to model building defines kNN as a lazy learner.
2. **Prediction at Test Time:** kNN's actual model building happens at prediction time. It analyzes the  $k$  nearest neighbors of the new data point and uses their class labels to predict the class for the new point. This "just-in-time" prediction process further solidifies kNN's lazy learner status.
3. **Memory Trade-off:** While kNN's training is faster than some eager learners due to the lack of upfront model building, it requires storing all the training data. This can lead to higher memory usage, especially for large datasets. This memory-prediction speed trade-off is another characteristic of lazy learners.

So, next time you encounter kNN, remember its delayed model building, just-in-time prediction, and memory-speed trade-off – all telltale signs of a true lazy learner!

## ⑤ Genetic algorithm.

— Bit string, what represents it?

→ for operation what?

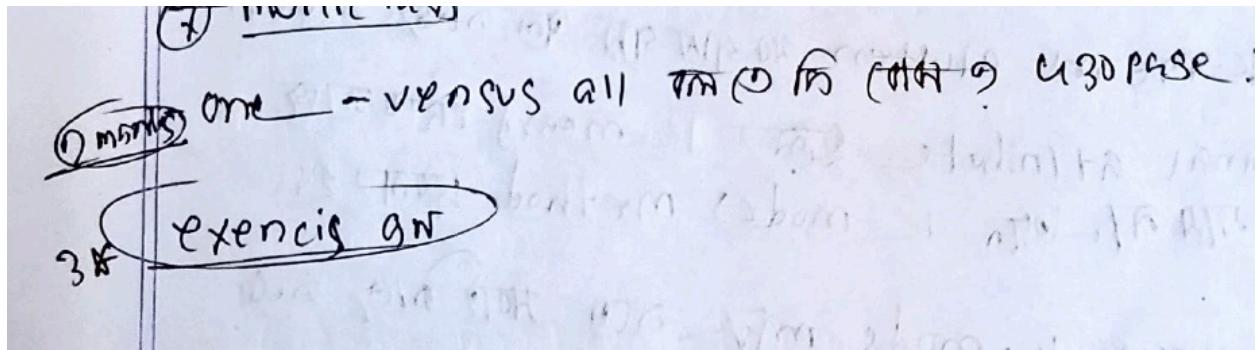
### 7.6.1 Genetic Algorithms

**Genetic algorithms** attempt to incorporate ideas of natural evolution. In general, genetic learning starts as follows. An initial **population** is created consisting of randomly generated rules. Each rule can be represented by a string of bits. As a simple example, suppose that samples in a given training set are described by two Boolean attributes,  $A_1$  and  $A_2$ , and that there are two classes,  $C_1$  and  $C_2$ . The rule “*IF  $A_1$  AND NOT  $A_2$  THEN  $C_2$* ” can be encoded as the bit string “100,” where the two leftmost bits represent attributes  $A_1$  and  $A_2$ , respectively, and the rightmost bit represents the class. Similarly, the rule “*IF NOT  $A_1$  AND NOT  $A_2$  THEN  $C_1$* ” can be encoded as “001.” If an attribute has  $k$  values, where  $k > 2$ , then  $k$  bits may be used to encode the attribute’s values. Classes can be encoded in a similar fashion.

Based on the notion of survival of the fittest, a new population is formed to consist of the *fittest* rules in the current population, as well as *offspring* of these rules. Typically, the **fitness** of a rule is assessed by its classification accuracy on a set of training samples.

Offspring are created by applying genetic operators such as crossover and mutation. In **crossover**, substrings from pairs of rules are swapped to form new pairs of rules. In **mutation**, randomly selected bits in a rule’s string are inverted.

The process of generating new populations based on prior populations of rules continues until a population,  $P$ , evolves where each rule in  $P$  satisfies a prespecified fitness threshold.



18. One vote. The class with the most votes is assigned.

**All-versus-all** (AVA) is an alternative approach that learns a classifier for each pair of classes. Given  $m$  classes, we construct  $\frac{m(m-1)}{2}$  binary classifiers. A classifier is trained

---

## 9.7 Additional Topics Regarding Classification 431

using tuples of the two classes it should discriminate. To classify an unknown tuple, each classifier votes. The tuple is assigned the class with the maximum number of votes. All-versus-all tends to be superior to one-versus-all.

A problem with the previous schemes is that binary classifiers are sensitive to errors. If any classifier makes an error, it can affect the vote count.