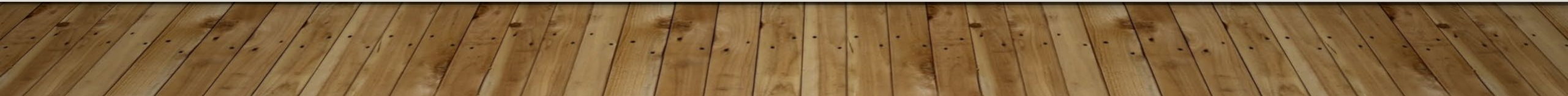


SUPERVISED LEARNING



MACHINE LEARNING LINGO

What?	Parameters	Structure	Hidden concepts	
What from?	Supervised	Unsupervised	Reinforcement	Self-supervised
What for?	Prediction	Diagnosis	Compression	Discovery
How?	Passive	Active	Online	Offline
Output?	Classification	Regression	Clustering	
Details??	Generative	Discriminative	Smoothing	

WHAT IS CLASSIFICATION?

A machine learning task that deals with identifying the class to which an instance belongs

A classifier performs classification



CLASSIFICATION LEARNING



Learning the classifier
from the available data
'Training set'
(Labeled)



Testing how well the classifier
performs
'Testing set'

AN EXAMPLE APPLICATION

- An emergency room in a hospital measures 17 variables (e.g., blood pressure, age, etc) of newly admitted patients.
- **A decision is needed:** whether to put a new patient in an intensive-care unit.
- Due to the high cost of ICU, those patients who may survive less than a month are given higher priority.
- **Problem:** to predict **high-risk patients** and discriminate them from **low-risk patients**.

ANOTHER APPLICATION

- A credit card company receives thousands of applications for new cards. Each application contains information about an applicant,
 - age
 - Marital status
 - annual salary
 - outstanding debts
 - credit rating
 - etc.
- **Problem:** to decide whether an application should approved, or to classify applications into two categories, **approved** and **not approved**.

MACHINE LEARNING AND OUR FOCUS

- Like human learning from past experiences.
- A computer does not have “experiences”.
- A computer system learns from data, which represent some “past experiences” of an application domain.
- Our focus: learn a target function that can be used to predict the values of a discrete class attribute, e.g., approve or not-approved, and high-risk or low risk.
- The task is commonly called: Supervised learning, classification, or inductive learning.

THE DATA AND THE GOAL

- **Data:** A set of data records (also called examples, instances or cases) described by
 - **k attributes:** A_1, A_2, \dots, A_k .
 - **a class:** Each example is labelled with a pre-defined class.
- **Goal:** To learn a **classification model** from the data that can be used to predict the classes of new (future, or test) cases/instances.

AN EXAMPLE: DATA (LOAN APPLICATION)

ID	Age	Has_Job	Own_House	Credit_Rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

Approved or not

AN EXAMPLE: THE LEARNING TASK

- Learn a classification model from the data
- Use the model to classify future loan applications into
 - Yes (approved) and
 - No (not approved)
- What is the class for following case/instance?

Age	Has_Job	Own_house	Credit-Rating	Class
young	false	false	good	?

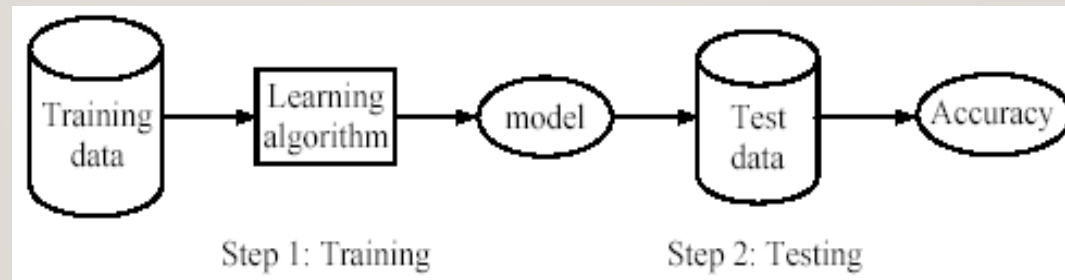
SUPERVISED VS. UNSUPERVISED LEARNING

- **Supervised learning:** classification is seen as supervised learning from examples.
 - **Supervision:** The data (observations, measurements, etc.) are labeled with pre-defined classes. It is like that a “teacher” gives the classes (**supervision**).
 - Test data are classified into these classes too.
- **Unsupervised learning (clustering)**
 - Class labels of the data are unknown
 - Given a set of data, the task is to establish the existence of classes or clusters in the data

SUPERVISED LEARNING PROCESS:TWO STEPS

- **Learning (training):** Learn a model using the **training data**
- **Testing:** Test the model using **unseen test data** to assess the model accuracy

$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}},$$



WHAT DO WE MEAN BY LEARNING?

- **Given**

- a data set D ,
- a task T , and
- a performance measure M ,

a computer system is said to **learn** from D to perform the task T if after learning the system's performance on T improves as measured by M .

- In other words, the learned model helps the system to perform T better as **compared to no learning**.

AN EXAMPLE

- **Data**: Loan application data
- **Task**: Predict whether a loan should be approved or not.
- **Performance measure**: accuracy.

No learning: classify all future applications (test data) to the majority class (i.e., **Yes**):

$$\text{Accuracy} = 9/15 = 60\%.$$

- We can do better than 60% with learning.

FUNDAMENTAL ASSUMPTION OF LEARNING

Assumption: The distribution of training examples is identical to the distribution of test examples (including future unseen examples).

- In practice, this assumption is often violated to certain degree.
- Strong violations will clearly result in poor classification accuracy.
- To achieve good accuracy on the test data, training examples must be sufficiently representative of the test data.

GENERATING DATASETS

- Methods:
 - Holdout (2/3rd training, 1/3rd testing)
 - Cross validation (n – fold)
 - Divide into n parts
 - Train on (n-1), test on last
 - Repeat for different combinations
 - Bootstrapping
 - Select random samples to form the training set

EVALUATING CLASSIFICATION METHODS

- **Predictive accuracy**

$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}}$$

- **Efficiency**

- time to construct the model
- time to use the model

- **Robustness**: handling noise and missing values

- **Scalability**: efficiency in disk-resident databases

- **Interpretability**:

- understandable and insight provided by the model

- **Compactness of the model**: size of the tree, or the number of rules.

EVALUATION METHODS

- **Holdout set:** The available data set D is divided into two disjoint subsets,
 - the *training set* D_{train} (for learning a model)
 - the *test set* D_{test} (for testing the model)
- **Important:** training set should not be used in testing and the test set should not be used in learning.
 - Unseen test set provides a unbiased estimate of accuracy.
- The test set is also called the **holdout set**. (the examples in the original data set D are all labeled with classes.)
- This method is mainly used when the data set D is large.

EVALUATION METHODS (CONT...)

- **n-fold cross-validation:** The available data is partitioned into n equal-size disjoint subsets.
- Use each subset as the test set and combine the rest $n-1$ subsets as the training set to learn a classifier.
- The procedure is run n times, which give n accuracies.
- The final estimated accuracy of learning is the average of the n accuracies.
- 10-fold and 5-fold cross-validations are commonly used.
- This method is used when the available data is not large.

EVALUATION METHODS (CONT...)

- **Leave-one-out cross-validation:** This method is used when the data set is very small.
- It is a special case of cross-validation
- Each fold of the cross validation has only a single test example and all the rest of the data is used in training.
- If the original data has m examples, this is m -fold cross-validation

EVALUATION METHODS (CONT...)

- **Validation set:** the available data is divided into three subsets,
 - a training set,
 - a validation set and
 - a test set.
- A validation set is used frequently for estimating parameters in learning algorithms.
- In such cases, the values that give the best accuracy on the validation set are used as the final parameter values.
- Cross-validation can be used for parameter estimating as well.

CLASSIFICATION MEASURES

- Accuracy is only one measure (error = 1-accuracy).
- **Accuracy is not suitable in some applications.**
- In text mining, we may only be interested in the documents of a particular topic, which are only a small portion of a big document collection.
- In classification involving skewed or highly imbalanced data, e.g., network intrusion and financial fraud detections, we are interested only in the minority class.
 - High accuracy does not mean any intrusion is detected.
 - E.g., 1% intrusion. Achieve 99% accuracy by doing nothing.
- The class of interest is commonly called the **positive class**, and the rest **negative classes**.

PRECISION AND RECALL MEASURES

- Used in information retrieval and text classification.
- We use a confusion matrix to introduce them.

	Classified Positive	Classified Negative
Actual Positive	TP	FN
Actual Negative	FP	TN

where

TP: the number of correct classifications of the positive examples (**true positive**),

FN: the number of incorrect classifications of positive examples (**false negative**),

FP: the number of incorrect classifications of negative examples (**false positive**), and

TN: the number of correct classifications of negative examples (**true negative**).

PRECISION AND RECALL MEASURES (CONT...)

	Classified Positive	Classified Negative
Actual Positive	TP	FN
Actual Negative	FP	TN

$$p = \frac{TP}{TP + FP} \quad r = \frac{TP}{TP + FN}$$

- **Precision** p is the number of **correctly classified positive examples** divided by the total number of examples that are classified as positive.
- **Recall** r is the number of **correctly classified positive examples** divided by the total number of actual positive examples in the test set.

AN EXAMPLE

	Classified Positive	Classified Negative
Actual Positive	1	99
Actual Negative	0	1000

- This confusion matrix gives

- precision $p = 100\%$ and
- recall $r = 1\%$

because we only classified one positive example correctly and no negative examples wrongly.

- **Note:** precision and recall only measure classification on the positive class.

F₁-VALUE (ALSO CALLED F₁-SCORE)

- It is hard to compare two classifiers using two measures. F₁ score combines precision and recall into one measure

$$F_1 = \frac{2pr}{p+r}$$

F₁-score is the harmonic mean of precision and recall.

$$F_1 = \frac{2}{\frac{1}{p} + \frac{1}{r}}$$

- The harmonic mean of two numbers tends to be closer to the smaller of the two.
- For F₁-value to be large, both p and r must be large.

RECEIVE OPERATING CHARACTERISTICS CURVE

- It is commonly called the **ROC curve**.
- It is a plot of the **true positive rate (TPR)** against the **false positive rate (FPR)**.
- True positive rate:

$$TPR = \frac{TP}{TP + FN}$$

- False positive rate:

$$FPR = \frac{FP}{TN + FP}$$

SENSITIVITY AND SPECIFICITY

- In statistics, there are two other evaluation measures:
 - **Sensitivity**: Same as TPR
 - **Specificity**: Also called **True Negative Rate** (TNR)

- Then we have

$$TNR = \frac{TN}{TN + FP}$$

$$FPR = 1 - specificity$$

EXAMPLE ROC CURVES

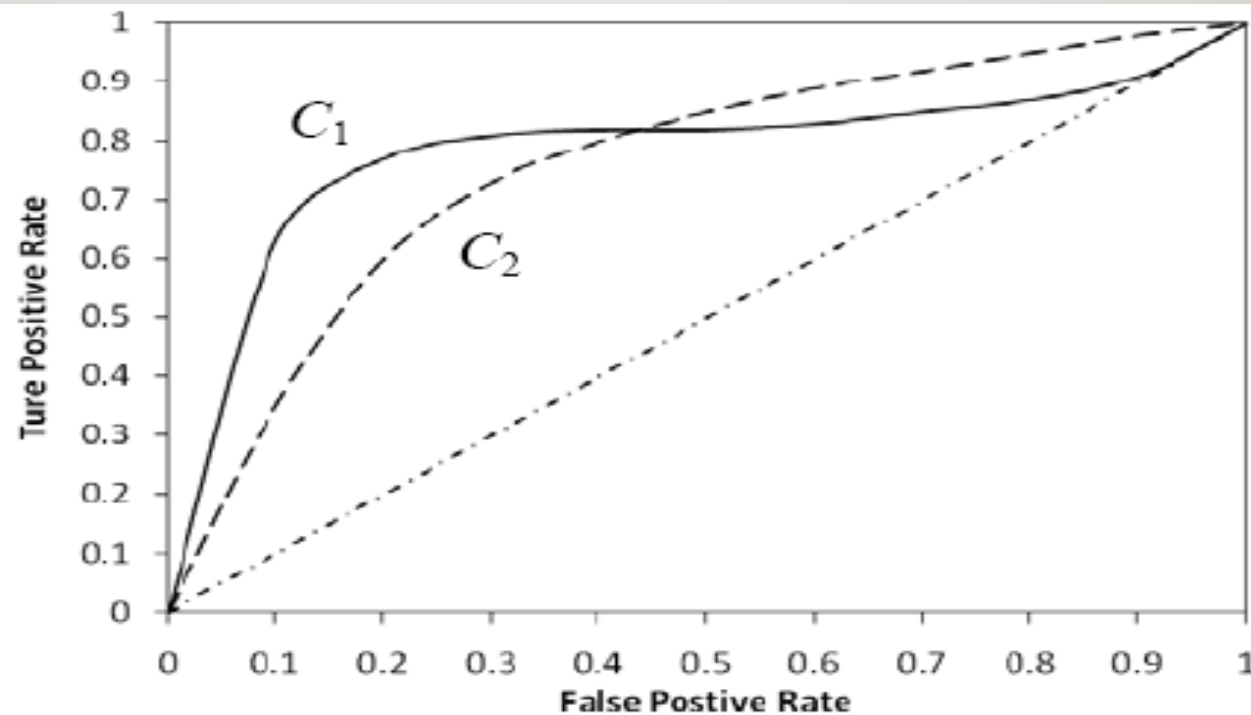


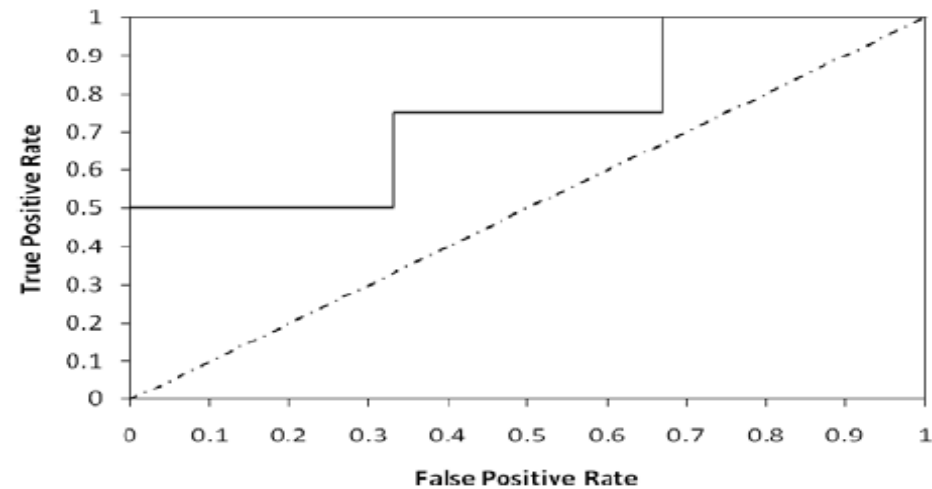
Fig. 3.8. ROC curves for two classifiers (C_1 and C_2) on the same data

AREA UNDER THE CURVE (AUC)

- Which classifier is better, C_1 or C_2 ?
 - It depends on which region you talk about.
- Can we have one measure?
 - Yes, we compute the area under the curve (AUC)
- If AUC for C_i is greater than that of C_j , it is said that C_i is better than C_j .
 - If a classifier is perfect, its AUC value is 1
 - If a classifier makes all random guesses, its AUC value is 0.5.

DRAWING AN ROC CURVE

Rank		1	2	3	4	5	6	7	8	9	10
Actual class		+	+	-	-	+	-	-	+	-	-
TP	0	1	2	2	2	3	3	3	4	4	4
FP	0	0	0	1	2	2	3	4	4	5	6
TN	6	6	6	5	4	4	3	2	2	1	0
FN	4	3	2	2	2	1	1	1	0	0	0
TPR	0	0.25	0.5	0.5	0.5	0.75	0.75	0.75	1	1	1
FPR	0	0	0	0.17	0.33	0.33	0.50	0.67	0.67	0.83	1



ANOTHER EVALUATION METHOD: SCORING AND RANKING

- **Scoring** is related to classification.
- We are interested in a single class (**positive class**), e.g., buyers class in a marketing database.
- Instead of assigning each test instance a definite class, scoring assigns a probability estimate (PE) to indicate the likelihood that the example belongs to the positive class.

RANKING AND LIFT ANALYSIS

- After each example is given a PE score, we can rank all examples according to their PEs.
- We then divide the data into n (say 10) bins. A lift curve can be drawn according to how many positive examples are in each bin. This is called **lift analysis**.
- Classification systems can be used for scoring. Need to produce a probability estimate.
 - E.g., in decision trees, we can use the confidence value at each leaf node as the score.

AN EXAMPLE

- We want to send promotion materials to potential customers to sell a watch.
- Each package cost \$0.50 to send (material and postage).
- If a watch is sold, we make \$5 profit.
- Suppose we have a large amount of past data for building a predictive/classification model. We also have a large list of potential customers.
- How many packages should we send and who should we send to?

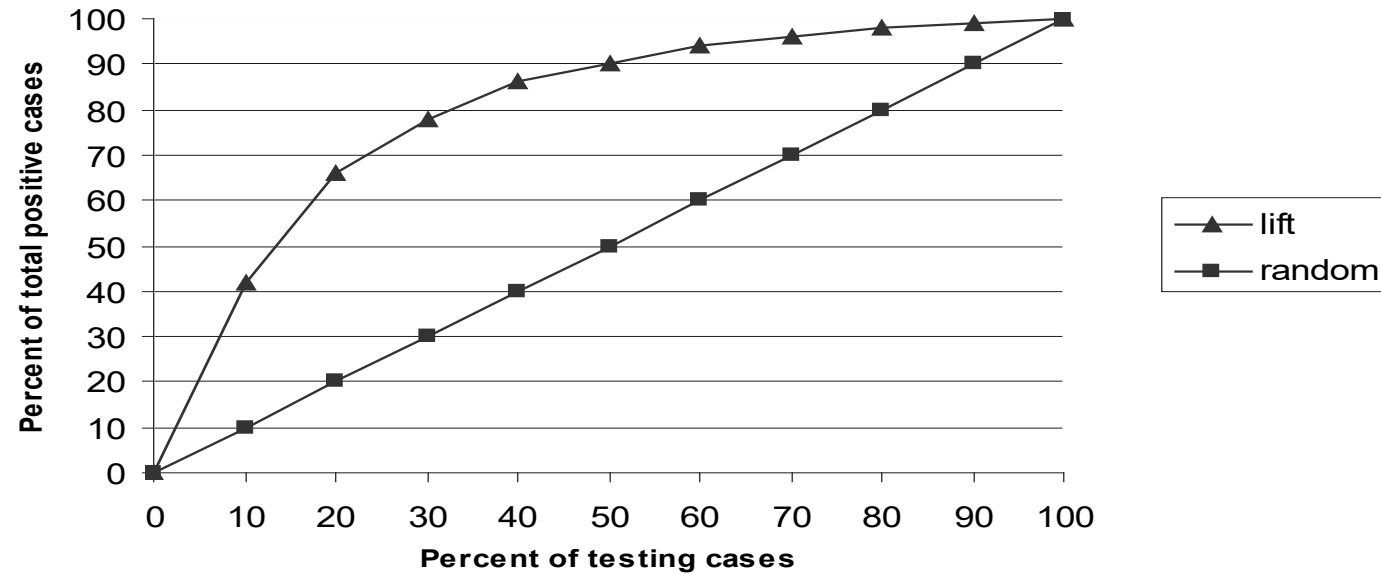
AN EXAMPLE

- Assume that the test set has 10000 instances. Out of this, 500 are positive cases.
- After the classifier is built, we score each test instance. We then rank the test set, and divide the ranked test set into 10 bins.
 - Each bin has 1000 test instances.
 - Bin 1 has 210 actual positive instances
 - Bin 2 has 120 actual positive instances
 - Bin 3 has 60 actual positive instances
 - ...
 - Bin 10 has 5 actual positive instances

LIFT CURVE

Bin	1	2	3	4	5	6	7	8	9	10
-----	---	---	---	---	---	---	---	---	---	----

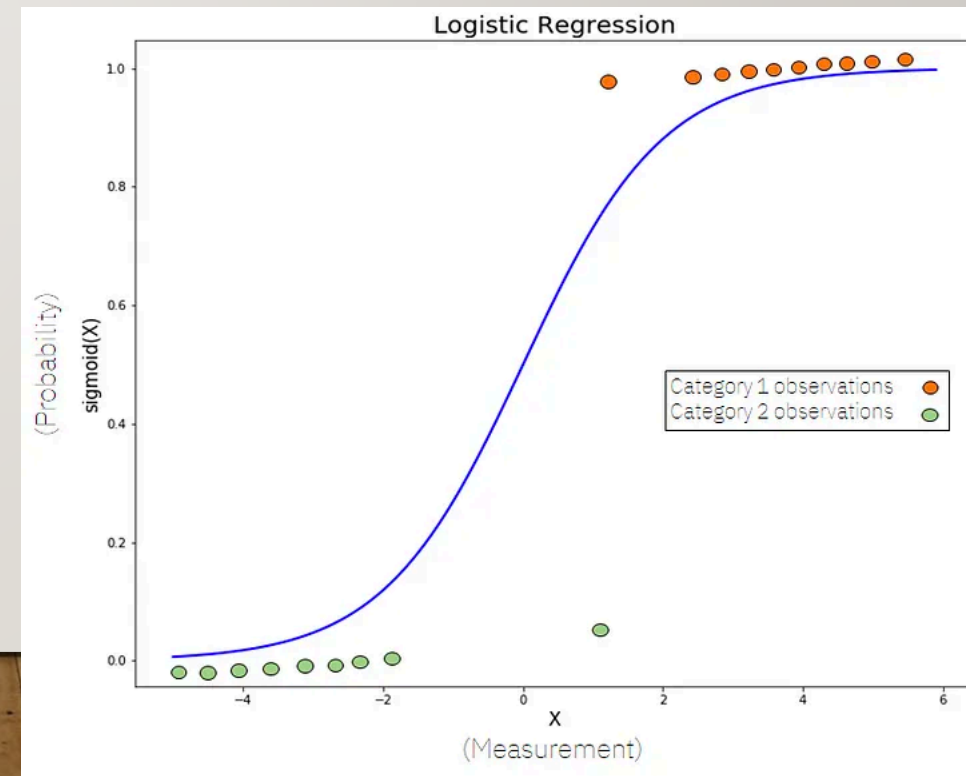
210	120	60	40	22	18	12	7	6	5
42%	24%	12%	8%	4.40%	3.60%	2.40%	1.40%	1.20%	1%
42%	66%	78%	86%	90.40%	94%	96.40%	97.80%	99%	100%



Eager Learners

CLASSIFICATION

- Classification puts records into groups.
- We examined logistic regression through the logit function.
- The model doesn't directly predict the label, it predicts the probability of that label, then translates that to an output.



Lazy learners

LAZY LEARNERS

- ‘**Lazy**’: Do not create a model of the training instances in advance
- When an instance arrives for testing, runs the algorithm to get the class prediction
- Example, K – nearest neighbour classifier
(K – NN classifier)
“One is known by the company
one keeps”

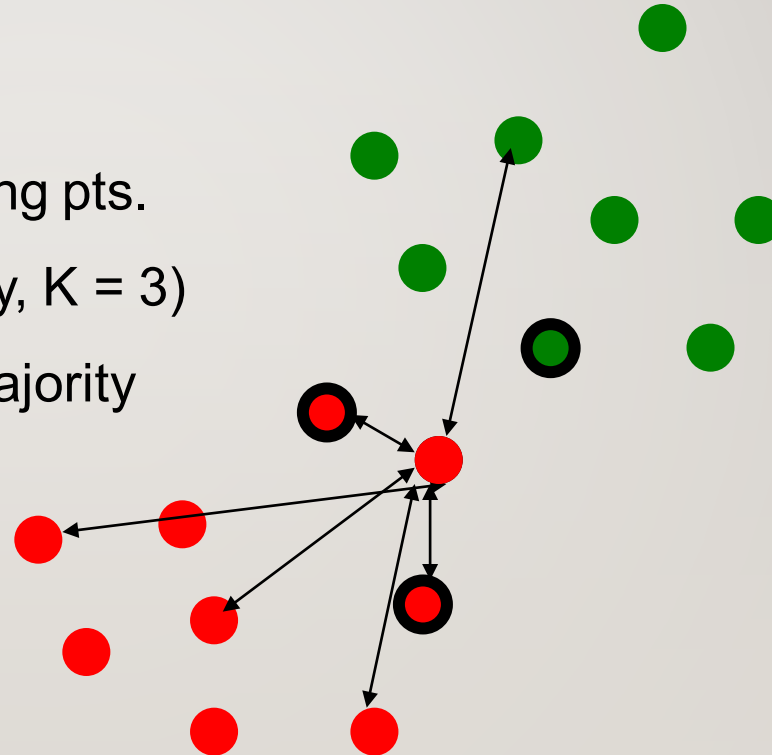
K-NN CLASSIFIER SCHEMATIC

For a test instance,

- 1) Calculate distances from training pts.
- 2) Find K-nearest neighbours (say, K = 3)
- 3) Assign class label based on majority

$$\text{dist}(X_1, X_2) = \sqrt{\sum_{i=1}^n (x_{1i} - x_{2i})^2}.$$

$$v' = \frac{v - \min_A}{\max_A - \min_A},$$



K-NN CLASSIFIER ISSUES

How to determine distances between values of categorical attributes?

Alternatives:

1. Boolean distance (1 if same, 0 if different)
2. Differential grading (e.g. weather – ‘drizzling’ and ‘rainy’ are closer than ‘rainy’ and ‘sunny’)