

Classification

Using sensor data to infer
semantic information

* Classification Problem

- Given a set of K classes

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$$

- and feature e

- learn a function f that assigns a class given feature.

$$C = f(e) \text{ with } C \in \Omega$$

* Learning Needs Information

- Learning the function f requires additional information.
- This additional information/knowledge can be provided through distributions.
- Manually specifying such knowledge is often hard and thus should be learned from data.

- Training data as a pairs of feature vector and classes.

$$(e, \omega) \xrightarrow{\text{Learning}} f$$

* Supervised Learning

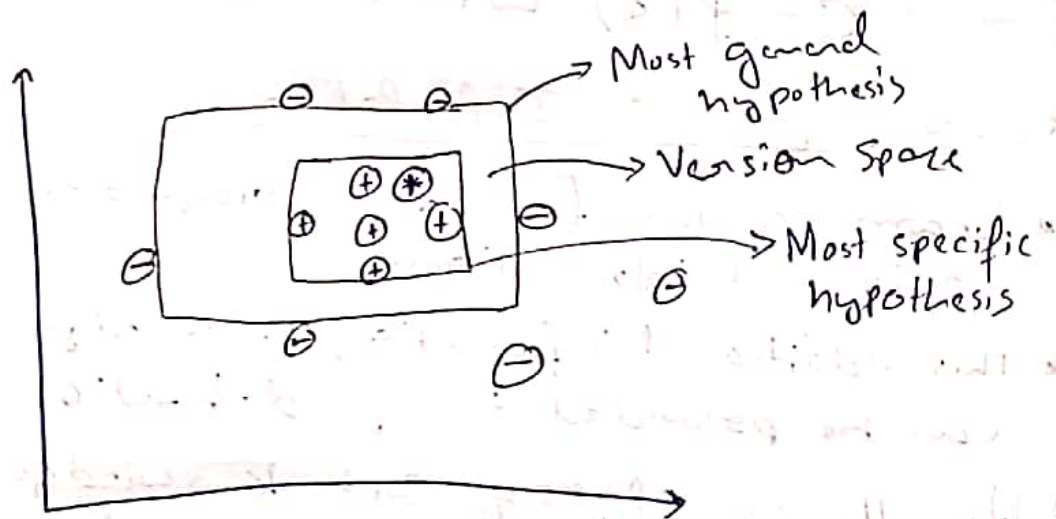
→ Learn a function f that generates / generalizes this decision to new (unseen) data.

* Classification

- Output is always a class label.
- Goal: Find a Separation of the input space that corresponds to the class.

* Regression

- Output is Continuous Variable
- Goal: Function that fits a set of data points.



⇒ Choose the hypothesis that maximizes the margin to the most specific and General one.

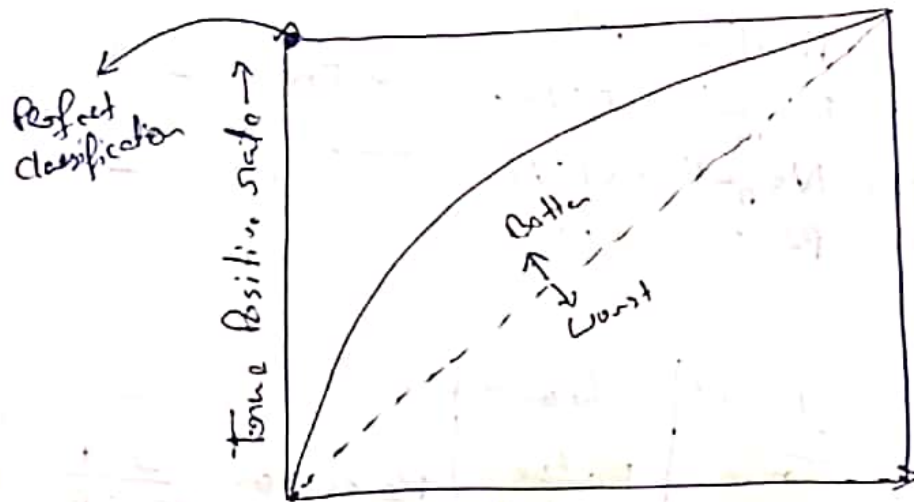
* Classification Error

⇒ Possible Outcomes:

- ① True Positive (TP)
- ② False Negative (FN) → Error I
- ③ True Negative (TN) → Error II
- ④ False Positive (FP)

	Condition +ve	Condition -ve	
Test outcome Positive	True Positive	False Positive	Precision = $\frac{\sum TP}{\sum OP}$
Test outcome negative	False Positive Negative	True Negative	Negative Prediction value = $\frac{\sum TN}{\sum ON}$
	Sensitivity = $\frac{\sum TP}{\sum CP}$	Specificity = $\frac{\sum TN}{\sum CN}$	Accuracy = $\frac{\sum TP + \sum TN}{\sum TP}$
	also called recall or true positive rate	also called true negative rate	

★ Receiver Operating Characteristic (ROC Curves)



False positive rate →



Classification

★ Designing a Classifier

Training Phase

Collecting labeled training data



Selecting and Computing appropriate features



Learning the distribution of features for each class or discriminant function.

Testing Phase (On different dataset)

Extract features
(as in the training phase)



Determine class based the classifier
(trained with a different dataset)



Evaluate the Performance

Operational Phase

Extract features



Determine class based the classifier

* Components of the Generalization Error

- * Bias: Describes how much the average model over all training set differ from the true model.
- * Variance: Describes how much models estimated from different training set differ from each other.

$$E(\text{MSE}) = \text{noise}^2 + \text{bias}^2 + \text{variance}$$

↓
{ Estimated mean square error }

Bias-Variance trade-off

* Underfitting

"Model is too 'simple' to represent the relevant characteristics"

- High bias and Low Variance
- High training error and high test error.

* Overfitting

Model is too "Complex" and fit irrelevant characteristics (noise) in the data.

- Low bias and high variance
- Low training error & high test error.

* Rules of Thumb

- Try simple classifiers first
- Use increasingly more powerful classifiers with more training data.

- Find good features

→ Better to have smart features and simple classifiers, than simple features and smart classifiers.

* 5x2 Cross Validation

- Randomly split up labeled dataset into 2 parts of equal size
- Use one for training, the second one for testing (validation)

- Swap both sets

- Repeat 5 times (called folds)

- Analyze the classification error

⇒ Results in different 10 classifiers.

* Nearest Neighbor Approach to Classification

⇒ The feature distribution is modeled by the training data (or a subset)

⇒ The class is assigned based on the closest feature in the training data.

Problem: Not Scalable

- Problematic in the presence of noise.
- Problematic for classes/features with different variances.
- Discriminant function is a Voronoi diagram



* NN and k-NN Approach

NN

→ The feature distribution is modeled by the training data (or a subset)

→ The class is assigned based on the closest feature in the training data.

K-NN

- The K nearest neighbors are considered for the classification decision (majority vote or weighted).
 - Choice of K is often done heuristically.
-

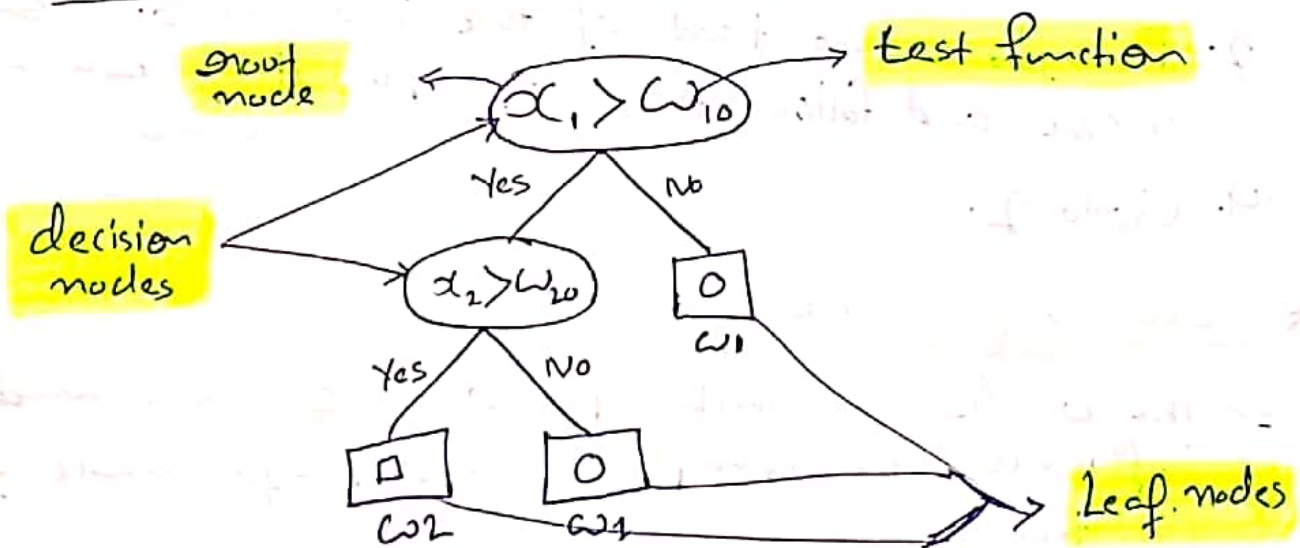
(11b)

Classification (cont)

* Decision Tree

- Idea \Rightarrow Sequence of splits of the input space define regions that correspond to classes.
- \Rightarrow Hierarchical data structure realizing a divide and conquer strategy.
- \Rightarrow Setup of the tree through training data.
- \Rightarrow Efficient nonparametric method for classification (and regression)

* Elements of a Decision Tree



* Decision Nodes

- \Rightarrow Each decision node implements a test function with discrete outcomes.
- \Rightarrow The test function of each decision node splits the input space into regions.
- \Rightarrow Also called split node.

* Leaf Node

- ⇒ A leaf node symbolizes the end of a sequence of decisions.
- ⇒ A single (output) class is associated to each leaf node.
- ⇒ A leaf node defines a localized region in the input space where instances falling in this region have the same label.

* Classification for a given decision Tree

1. Start at the root node.
2. If current node is a leaf node, return its class label.
3. Perform the test of the current decision node and follow the corresponding branch.
4. Goto 2.

* Learning a Decision Tree

- ⇒ The order in which split decisions are made influences the complexity and performance of the tree.
- ⇒ Finding the optimal arrangement of tests is NP hard, thus heuristics are needed.

* Which Decision to Make Next

- ⇒ Select the test that best Separates the Class labels in the data.
- ⇒ The "Purer" the children, the better the split.
- ⇒ For any "pure" child branch, we can create a leaf node (no further splits).

* Impurity

- ⇒ Purity (or impurity) can be defined through the uncertainty in the distribution over the class label in the current vs the split-up region.
- ⇒ Goal: Always select the split that minimizes impurity.
- ⇒ Different impurity measures:
 - Entropy
 - Gini index

* Entropy as Impurity

→ Measure of uncertainty in probability distribution.

$$H(\Omega) = - \sum_{\omega_i \in \Omega} P(\omega_i) \ln P(\omega_i)$$

- ⇒ Select the split that reduces the entropy of most:

$$\Delta H = H(\Omega) - (P_{\text{left}} H_{\text{left}}(\Omega) + (1 - P_{\text{left}}) H_{\text{right}}(\Omega))$$

change in
entropy

Entropy
before
split

Entropy after split.

* Gini Index

→ Alternative Criterion to entropy.

→ Gini index

$$G(\Omega) = 1 - \sum_{c_i \in \Omega} p(c_i)^2$$

* When to Stop?

▪ Intuitive idea: add a leaf node of a split leads to a pure node.

▪ Overfitting problem: The tree perfectly splits the classes on the training dataset but does not generalize well to other dataset.

⇒ Standard approach: Stop after a certain level of purity is reached.

⇒ A leaf stores the posterior probabilities of classes, instead of best label.

* Decision Tree for Classification

→ Comparably easy to understand and implement.

→ Works well to high-dimensional data.

→ Allows to handle numerical and categorical variables easily.

→ Finding the optimal split is NP hard.

→ Heuristics are used (e.g., entropy)

* Maximum-A Posteriori Classification

- Classification is decision making
- Probability theory as the framework for making decisions under uncertainty.
- Based on Bayes' rule.

$$P(w|e) = \frac{P(e|w) P(w)}{P(e)}$$

Class feature

(a-priori Probability for the occurrence of the class)

"What is the probability of a class given an observed feature"

(Probability for the feature occurrence can be obtained by:

(Likelihood function for the class w . It is also called observation model)

$$P(e) = \sum_{k=1}^K P(e|w_k) P(w_k)$$

$$\text{Posterior} = \frac{\text{likelihood} \times \text{Prior}}{\text{evidence}}$$

$$= \eta \times \text{likelihood} \times \text{Prior}$$

⇒ Distributions $P(e|\omega)$ and $P(\omega)$ must be learned from training data.

* Map Classification

1. Compute for each class

$$P(\omega_i | e) = \frac{P(e | \omega_i) P(\omega_i)}{\sum_{k=1}^K P(e | \omega_k) P(\omega_k)}$$

2. Select the most likely class

$$\omega_{i^*} \text{ with } i^* = \arg \max_i P(\omega_i | e)$$

* MAP Classification with Gaussian Distributed Features

⇒ Let us look into features that follow a Gaussian given a class

$$e = x \quad x | \omega_i \sim \mathcal{G}(\mu_i, \Sigma_i) \quad i = 1, 2$$

⇒ Thus, we can write

$$P(\omega_i | x_i) = \eta \mathcal{G}(x, \mu_i, \Sigma_i) P(\omega_i)$$

⇒ and the negative log likelihood

$$-\ln P(\omega_i | x) = -\ln \eta - \ln \mathcal{G}(x, \mu_i, \Sigma_i) - \ln P(\omega_i)$$

⇒ The MAP classifier directly yields

$$\arg \max_{\omega_i} P(\omega_i | x)$$

$$\Rightarrow \arg \min_{\omega_i} \left(-\ln g(x, \mu_i, \Sigma_i) \right) - \ln P(\omega_i)$$

⇒ and the classification boundary are points in which the function

$$D(x) = -\ln(P(\omega_1 | x)) + \ln(P(\omega_2 | x))$$

changes its sign

⇒ The shape of $D(x) = 0$ depends on the eigenvalues of $\Sigma_1^{-1}, \Sigma_2^{-1}$.

~~* Designing a Classifier: Training Phase~~

* Learning a Classifier

→ In practice the distributions $P(x | \omega)$ and $P(\omega)$ are not known.

→ Both must be learned from data

* Training data

▪ Sample set of size N

▪ N_i Sample corresponds to class ω_i

* Prior Distribution ($P(\omega_i)$)

⇒ The prior distribution $P(\omega_i)$, which models the probability that a random sample corresponds to class ω_i :

$$P(\omega_i) = \frac{N_i}{N}$$

* The Likelihood Function

$$P(x_i | \omega_i) = g(x_i, \mu_i, \Sigma_i)$$

$$\mu_i = \frac{1}{N_i} \sum_{j=1}^{N_i} x_{ij}$$

$$\Sigma_i = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (x_{ij} - \mu_i)(x_{ij} - \mu_i)^T$$

Generative Approach

→ Use data to calculate the posterior densities that get the discriminant function.

Discriminative Approach

→ Bypasses the estimation of densities & directly estimate the discriminants.

