supervised learning: learning model from label
unsupervised learning: extract meaningful information without label
(clustering / dimensionality reduction)
Reinforced learning: improve the performance based on interaction with environment. The feedback is a measure of how well the action was measured by a reward function.

| | | |
|---|---|---|
| Decision Tree | non | |
| Linear SVM | linear | |
| SVM with kernel | non | |
| Linear Regression | a classifier | |
| Logistic Regression | linear | |
| Naive Bayes | linear | |

## Linear Regression : Least Square Fitting

minimize : $\sum_{i}^{n} [y_i - (\beta_0 + \beta_1 x_i)]^2$

$S_{xy} = \sum_{i=1}^{n} x_i y_i - \frac{1}{n} (\sum_{i=1}^{n} x_i)(\sum_{i=1}^{n} y_i)$

$S_{xx} = \sum_{i=1}^{n} x_i^2 - \frac{1}{n} (\sum_{i=1}^{n} x_i)^2$

$S_{yy} = \sum_{i=1}^{n} y_i^2 - \frac{1}{n} (\sum_{i=1}^{n} y_i)^2$

$\hat{\beta_0} = \bar{y} - \hat{\beta_1} \bar{x}$       $\hat{\beta_1} = \frac{S_{xy}}{S_{xx}}$

Multiple: $\hat{\beta} = (X^T X)^{-1} X^T Y$

## Decision Trees:

when have $N$ attributes. $2^{(2^n)}$ trees

Choose feature $X_i = H(Y|x_i) \downarrow$
   $I(X;Y) = H(Y) - H(Y|X) \uparrow$

$H(x) = -\sum p(x) \log_2 p(x)$

$H(Y|x) = \sum p(x) H(Y|X=x)$

$I(x,Y) = \sum P_{x,y} (x,y) \log \left( \frac{P_{x,y} (x,y)}{P_x(x) P_y(y)} \right)$

$Gini = 1 - \sum P_j^2$

classification Error = $1 - \max P_j$

Decision Tree Depth $\uparrow \to$ overfitting

## Ensemble learning - Random Forest
reduce overfitting and variance without decreasing performance

## Bagging - Bootstrap Aggregating

Bootstrapping: Random sampling with replacement

train multiple decision trees & search all features to split on for each tree

Aggregating: Combine multiple predictions via averaging or majority vote

---

## Support Vector Machines

kernel : mapping to higher-dimensional space

complexity depends on the number of training samples, not dimensionality

larger margin → lower error

$\phi(w) = w^T w$ is minimized

Soft margin is allowed errors.

$\phi(w) = w^T w + \boxed{C} \sum \xi_i$

penalty $C \to 0$  underfitting
       $C \to \infty$  overfitting

### SVM Weakness :
① sensitive to noise
② Standard SVM only consider 2 classes → build multiple SVMs;
③ select a specific kernel and parameters is usually done by see and try

## Naive Bayes → MAP

$\underset{y}{\arg\max} P(y|x) = \underset{y}{\arg\max} P(y) \prod_{i=1}^{n} P(x_i|y)$

{ Logistic Regression directly compute $P(y|x)$
  Naive Bayes use Bayes Theorem to compute $P(y|x)$

NB + Gaussian Basis Function
× LR + sigmoid

Logistic Regression model the $P(y|x)$ as a logistic function. The logic is a weighted linear combination of the features. Linear regression itself is based on linear feature function to regress.

→ Or using total probability:

$P(S) = P(S|yes) P(yes) + P(S|no) P(no) = (2/9)(9/14) + (3/5)(5/14) = \frac{5}{14}$

NB: $P(yes|S,w) = P(S,w|yes) P(yes)/P(S,w) = \frac{10}{37}$

NB: $P(S,w|yes) = P(S|yes) P(w|yes) = \frac{6}{81}$

→ $P(S,w) = P(S,w|yes) + P(S,w|no) P(no)$
       $= P(S|yes) P(w|yes) P(yes) + P(S|no) P(w|no) P(no)$
       $= \frac{37}{240}$

NB: input features $X_i$ is independent given label $Y$

LR: Least square is not suitable. Maximum likelihood Estimation instead

---

## Performance Matrix

Actual value

| | | |
|---|---|---|
| Predicted value | TP | FP |
| | FN | TN |

(TPR) Recall : $\frac{TP}{FN+TP}$  when false negatives is catastrophic. e.g. disease detection

(PPV) Precision : $\frac{TP}{TP+FP}$  when being right (positive prediction is correct) outweighs detecting all positives. e.g. recommendation system

$NPV = \frac{TN}{TN+FN}$

① Training set:
↓
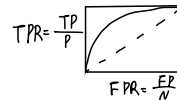② Validation set: hyperparameters tuning and model selection
↓
③ Test data

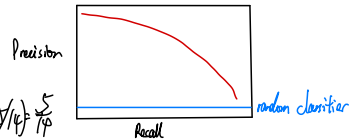K-fold cross validation is used when we have little data



F1- score: $2 \times \frac{Precision \times Recall}{Precision + Recall}$

(TNR) Specificity = $\frac{TN}{TN + FP}$

ROC Curves : the trade-off between TP rate and FP rate



$TPR = \frac{TP}{P}$

$FPR = \frac{FP}{N}$

PR Curve: between TP rate and positive predictive value



Precision / Recall / random classifier

Occam's Razor
Given 2 models with similar generalization errors, the simple model is preferred.

Probabilistic classifiers : yields a probability distribution. includes: Random forest, NB, LR

Deterministic classifiers : no model → separate feature space. includes: Decision Tree, SVM

Generative classifiers : learn $P(x,y)$, calculate $P(y|x)$ to find arg max $P(y)$ NB, Bayesian Network, Hidden Markov

Discriminative classifiers : learn $P(y|x)$ directly LR, SVM

*NN*: supervised learning

perceptron : $y = 6\left(b + \sum_{i=1}^{m} w_i x_i\right) = \{0, 1\}$

Multilayer perceptron: $y_j = 6\left(\sum_{i=1}^{m} w_{ij} x_i + b_j\right)$

$\hookrightarrow$ feed forward network

Common Activation Functions:

  Step / Sigmoid / tanh / ReLu / Leaky ReLu

  Maxout / ELU

NN = Function Approximation

Feedforward NN : no loops   input → hidden layers → output

Recurrent : use feedback

RNN: ① employ feedback
        not necessarily stable

    ② forecasting time series data
       language translation

Hopfield Networks: fully connected

 CNN → Image

Transformer Networks : Natural language Processing

Generative Adversarial Network:

---

Single perceptron can find linear max margin ✓
            LR solution ✓

---

$NB \rightarrow$ interest in $P(labels | features)$

$B: P(S | (F_1, F_2) = P((F_1, F_2) | S) * P(S) / P(F_1, F_2)$

$NB: P(S | (F_1, F_2) = P((F_1, F_2) | S) \times P(S) / P(F_1, F_2)$

$$= \frac{P(F_1 | S) \times P(F_2 | S) \times P(S)}{P(F_1) \times P(F_2)}$$

*NN* optimization: ① Gradient Descent
                      ② Backpropagation

Computing & Intelligence drive IoT

IoT Architecture:

Application Domain ⟷ Network D ⟷ Device D

{ apps          { M2M                    { Device
  ...           Communication Network     M2M Area Network

            Access Network          IoT Gateway

Bottom End : Basic. Resource constrained, environment monitor

In the Middle: support localization. full protocol. in Home, industrial

Top End : Military / Medical USE

Technical Issue: ① Naming, Addressing, Routing
              ② Scalability for network architecture
              ③ Power saving & management
              ④ Security

Challenges : ① Stringent Latency Requirements
          ② Network Bandwidth Constrain
          ③ Resource Constrained Device