**Decision Tree:**

$$H(v) = \sum_V -P(H=v)\log_2 P(H=v).$$

$$Gain(S,A) = H(S) - \sum_{v \in Value(A)} P(v) \cdot H(v)$$

**Point Estimation.**

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}. \qquad \hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

**Naive Bayes:**

$$P(Y=y_{v} | X=X_R) = \frac{P(X=X_R | Y=y_v) \cdot P(Y=y_v)}{\sum_j P(X=X_R | Y=y_v) \cdot P(Y=y_v)}$$

**Logistic Regression.**

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_v w_v x_v)}.$$

**Perceptron. Training Rule:**

$$\begin{cases} \Delta w_i = \eta(t-o) \cdot x_i \\ w_i = w_i + \Delta w_i. \end{cases}$$

To guarantee convergence; use gradient descent:

minimize $E[\vec{W}] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$

$$\nabla E[\vec{W}] = [\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n}].$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \cdot \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$= \sum_{d \in D} (t_d - o_d)(-x_{i,d})$$

$$\vec{\Delta w} = -\eta \cdot \nabla E[\vec{W}]. \text{ therefore,}$$

$$\Delta w_i = -\eta \cdot \sum_{d \in D} (t_d - o_d)(-x_{i,d})$$

$$= \eta \cdot \sum_{d \in D} (t_d - o_d) \cdot x_{i,d}.$$

**Batch_Gradient_descent(D,η):**

1. $\vec{W} = <w_0, w_1, \dots, w_n> \leftarrow$ small random number.
2. Until convergence do.
   - i). $\vec{\Delta w_i} = <\Delta w_0, \dots, \Delta w_n> \leftarrow 0$
   - ii). for $d \in D$, where $\vec{d} = <\vec{x}, y>$, do
     $$o_d = w^T \cdot \vec{x}$$

---

for $\Delta w_i \in \Delta \vec{W_i}$, do
$$\Delta w_i \leftarrow \Delta w_i + \eta(y - o_d) \cdot x_i$$

iii). for $w_i \in \vec{w}$, do
$$w_i \leftarrow w_i + \Delta w_i$$

3. return $\vec{w}$.

**Incremental_Gradient_Descent(D,η)**

1. $\vec{W} = <w_0, w_1, \dots, w_n> \leftarrow$ small random value.
2. Until meet the termination condition do
   for $d \in D$ do
   i) $o_d = \vec{w}^T \cdot \vec{x}$.
   for $\Delta w_i \in \Delta \vec{W_i}$ do
   $$\Delta w_i = \eta \cdot x_i \cdot (y - o_d)$$
   for $w_i \in \vec{w}$, do
   $$w_i = w_i + \Delta w_i$$

3. Return $\vec{w}$.

---

If $o_d = $ Sigmoid$(\vec{w} \cdot \vec{x}) = \vec{\sigma}(\vec{w} \cdot \vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$

then $\frac{d \vec{\sigma}(z)}{dz} = \vec{\sigma}(z) \cdot (1 - \vec{\sigma}(z))$

Then $\Delta w_i = -\eta \cdot \frac{\partial E}{\partial w_i} = \eta (t_d - o_d) \cdot o_d \cdot (1 - o_d) \cdot x_i.$

**Neural Network.**

$$o_v \xrightarrow{w_{ji}} o_{j/h} \xrightarrow{w_{kh}} o_k.$$

**Backpropagation(D, η, $n_{in}$, $n_{out}$, $n_{hidden}$)**

1. Create a feed-forward network
2. initialize all network weights to small random number
3. Until the termination condition is met, do
   for $d \in D$ do
   i) compute output $o_u$ of every unit.
   ii). for each output unit k, its error term
   $$\delta_k = \underset{f'}{\underline{o_k(1-o_k)}}(t_k - o_k)$$
   iii). for each hidden unit h,
   $$\delta_h = \underset{f'}{\underline{o_h(1-o_h)}} \sum_{k \in outputs} w_{k,h} \cdot \delta_k$$
   iii)
   $$w_{ji} = w_{ji} + \eta \cdot \delta_j \cdot x_{ji}$$

# Instance Based Learning.

N points in D dimensions, to apply k-nearest neighbor algorithm. the distance must go is: $\left(\frac{k}{N}\right)^{\frac{1}{D}}$

## SVM:

Lagrangian function.

i). Primal form: minimize

$$L_p(\vec{w}, b, \alpha_i) = \frac{1}{2}||w||^2 - \sum_{v=1}^{n} \alpha_i \cdot (y_i \cdot (\vec{w}_i \cdot x_i + b) - 1)$$

such that $\alpha_i \geq 0$.

KKT conditions:

only support vectors have $\alpha_i \neq 0$, then get b from ③

① $\frac{\partial L_p}{\partial w} = 0 \Rightarrow \vec{w} = \sum_{v=1}^{n} \alpha_i \cdot y_i \cdot \vec{x_i}$

② $\frac{\partial L_p}{\partial b} = 0 \Rightarrow \sum_{v=1}^{n} \alpha_i \cdot y_i = 0$

③ $\alpha_i \cdot (y_i \cdot (\vec{w} \cdot \vec{x_i} + b) - 1) = 0$.

④ $\alpha_i \geq 0$.

ii). Lagrangian Dual Problem.

maximize $\sum_{v=1}^{n} \alpha_i - \frac{1}{2} \sum_{v=1}^{n} \sum_{j=1}^{n} \alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot \vec{x_i}^T \cdot \vec{x_j}$

such that $\alpha_i \geq 0$, $\sum_{v=1}^{n} \alpha_i \cdot y_i = 0$.

Kernel function $K(x_i, x_j) = ?$

After solve $\alpha_i$, given a new test data $x$,

the prediction $g(x) = \sum_{i \in SV} \alpha_i \cdot K(x_i, x) + b = \begin{cases} 1, & + \\ -1, & - \end{cases}$

---

## Ensemble methods. (Bagging and boosting).

AdaBoost:

1. initialize the data $\{W_n\}$ : $W_n^{(1)} = \frac{1}{N}$.

2. For $m = 1, \cdots, M$.

   i). minimize the weighted error function.

   $$J_m = \sum_{n=1}^{N} W_n^{(m)} \cdot I(y_m(x_n) \neq t_n)$$

   ii).
   
   $$\varepsilon_m = \frac{J_m}{\sum_{n=1}^{N} W_n^{(m)}}, \quad \alpha_m = \ln\left\{\frac{1 - \varepsilon_m}{\varepsilon_m}\right\}$$

   iii). $W_n^{(m+1)} = W_n^{(m)} \cdot e^{\{\alpha_m \cdot I(y_m(x_n) \neq t_n)\}}$

3. $Y_M(x) = sign\left(\sum_{m}^{M} \alpha_m \cdot y_m(x)\right)$

---

# Bayesian Network.

Variable elimination: process nodes in order, then eliminate = connect all children of a node to each other.

Learning theory:

$$m \geq \frac{1}{\varepsilon}\left(\ln\left(\frac{1}{\delta}\right) + \ln|H|\right)$$

$$m \geq \frac{1}{2\varepsilon^2}\left(\ln\left(\frac{1}{\delta}\right) + \ln|H|\right)$$

$$m \geq \frac{1}{\varepsilon}\left(4\log_2\left(\frac{2}{\delta}\right) + 8 VC(H) \cdot \log_2\left(\frac{13}{\varepsilon}\right)\right)$$

---

The VC-dimension is always less than the size of hypothesis since $|H| \geq 2^m$, to shatter m instances, $VC(H) \leq \log_2(H)$.

At each iteration, EM improves its cost function and is guaranteed to converge to a local minima.

Gaussian Naive Bayes is a linear classifier only when the variances are assumed to be independent of the class.

No free lunch theorem.

Gaussian naive Bayes classifier has more expressive power than a linear classifier if we don't restrict the variance parameters to be class independent.

Add the constraint that $0 < \alpha_i \leq c$, SVM.

Slower than decision tree: Bagging, boosting.

Bagging can be parallelized easily, but hurt the performance of stable class

Bagging and boosting attempt to reduce variance.

F. if the error $\varepsilon > 0.5$, then the weight assigned to the misclassify points will be smaller than the weight assigned to the points that are classified correctly, then the subsequent iterations will not try to classify these points correctly and thus the algorithm is likely to exhibit poor performance.

Adaboost is susceptible to outliers, a possible heuristic is to put a threshold on the weight and remove all points that have very large weights.

KNN: as we increase k, the algorithm underfits because the decision surface becomes simpler.

Perceptron training rule: (also for linear regression,

$$\begin{cases} \Delta W_i = \dfrac{\partial E}{\partial W_i} \quad , \text{ where } E(w) = \dfrac{1}{2}\sum_{d \in D}(t_d - o_d)^2 \\ W_i \leftarrow W_i + \eta \cdot \Delta W_i \end{cases}$$

① Batch gradient descent algorithm:

    initialize $\vec{w} = <W_0, W_1, \cdots, W_n> \leftarrow$ small random value

    Until convergence, do

        initialize $\Delta \vec{W} = <\Delta W_0, \Delta W_1, \cdots, \Delta W_n> \leftarrow 0$

        for $d \in D$, do

            $o_d = \vec{w}^T \cdot \vec{x}$.

            for $i = 0$ to $n$, do

                $\Delta W_i \leftarrow \Delta W_i + \dfrac{\partial E}{\partial W_i}$

        for $i = 0$ to $n$, do

            $W_i \leftarrow W_i + \eta \cdot \Delta W_i$

    return $\vec{w}$.

Support Vector Machine.

    Lagrangian dual problem: maximize $\sum_{i=1}^{n}\alpha_i - \dfrac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot x_i^T \cdot x_j$

      with slack penalty.

                    such that $C \geq \alpha_i \geq 0$, and $\sum_{i=1}^{n}\alpha_i \cdot y_i = 0$

$$L_p(\vec{w}, b, \alpha_i) = \frac{1}{2}\|w\|^2 - \sum_{i=1}^{n}\alpha_i \cdot (y_i \cdot (\vec{w}^T \cdot \vec{x}_i + b) - 1), \ \alpha_i \geq 0$$

$$g(x) = \vec{w}^T \cdot x + b = \sum_{i \in SV}\alpha_i \cdot x_i^T \cdot x + b$$

all $\alpha_i$ should be bounded, if the optimal $\alpha$ are greater than $C$, we may end up with a sub-optimal decision boundary.

As data points $N \to \infty$, prior $P(\theta)$ is forgotten, and MLE = MAP.

Point Estimator:

    $\begin{cases} \text{(i). MLE: } P(D|\theta) \quad , \ \hat{\theta} = \arg\max_\theta \ln P(D|\theta) = \dfrac{d}{d\theta}\ln P(D|\theta) = 0. \\[4mm] \text{(ii). MAP: } P(\theta|D) \quad , \ \hat{\theta} = \arg\max_\theta P(\theta|D) = \dfrac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2} \end{cases}$

    We use MAP in naive Bayes.

Naïve Bayes: $P(Y = y_i | X = x_k) = \dfrac{P(X = x_k | Y = y_i) \cdot P(Y = y_i)}{\sum_j P(X = x_k | Y = y_j) \cdot P(Y = y_j)}$

if we assume the variance is independent of the class, then naive Bayes is a linear classifier.

$$= \dfrac{P(X_1 \cdots X_n | Y = y_i) \cdot P(Y = y_i)}{\sum_j \prod_i P(X_i | Y = y_i) \cdot P(Y = y_i)}$$

$\text{Beta}(\beta_H, \beta_T)$

① Gaussian naive Bayes: estimate the parameters: $P(Y = y_k)$, $u_{ik} = E[X_i | Y = y_k]$, $\sigma_{ik}^2 = E[(X_i - u_{ik})^2 | y_k]$

② Discrete-valued Input: $P(X_i | Y = y_k)$, $P(Y = y_k)$.

---

Decision Tree:

$H(v) = \sum -P(v) \cdot \log P(v)$

$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{Value}(A)} P(v) \cdot H(S_v)$

Given $X = (X_1, \cdots, X_d)$, the bound on the number of leaf nodes is $O(2^d)$.

$\dfrac{d\sigma(x)}{dx} = \sigma(x) \cdot (1 - \sigma(x))$

# Neural Network·Backpropagation :

↑ the error surface may contain many different local minima, therefore, only guaranteed to converge toward some local minima, not global minima.

$$E_d(\vec{w}) = \frac{1}{2}|t_d - O_d|^2$$

1. initialize all network weights to small random number.

2. Until convergence do

    for each $<x, y>$ in training-example do.

        i). input the instance $\bar{x}$ to the network, and compute the output $O_u$ of every unit u in the network.

        ii). for each network |output| unit k, calculate its error term $\delta_k$.

$$\delta_k \leftarrow O_k(1-O_k)(t_k-O_k) \qquad // \quad \frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j}$$

$$\delta_k = -\frac{\partial E_d}{\partial net_k}$$

        iii). for each hidden unit h, calculate its error term $\delta_h$.

$$\delta_h \leftarrow O_h(1-O_h) \sum_{k \in outputs} W_{k\cdot h}\delta_k$$

output of h.

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in outputs} \left(\frac{\partial E_d}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j}\right)$$

$$= \sum_{k \in outputs} -\delta_k \cdot \frac{\partial net_k}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j}$$

$$= \sum_{k \in outputs} -O_j(1-O_j) \cdot W_{k\cdot j} \cdot \delta_k$$

        iv). Update each network weight $w_i$

$$W_i \leftarrow W_i + \eta \cdot \delta_j \cdot x_{ji}$$

    ↓ input
    target

$$\Delta w_{ji} = -\eta \cdot \frac{\partial E_d}{\partial w_{ji}}$$

$x_0$

$\delta$   $\downarrow w_1$

$$w_1 \leftarrow w_1 + \eta \cdot \delta \cdot x_0$$

---

# Logistic Regression MCAP. L2 algorithm:

the objective function is always concave, and therefore the algorithm is always converge under mild assumption.

    initialize $w = <w_0, w_1, \cdots, w_n>$.

    until convergence do

        i) for each data $d \in D$.

$$prob[d] = sigmoid(W^T \cdot x)$$

        ii). initialize $dw = <dw_0, dw_1, \cdots, dw_n>$ to 0.

        iii). for i=0 to n

            for j=0 to m

$$dw[i] = dw[i] + x_{ji} \cdot (Y_j - Prob[j]) = \frac{\partial Loss(w)}{\partial w_i}$$

        iv). for i=0 to n

$$w[i] = w[i] + \eta(dw[i] - \lambda \cdot w[i])$$

The time complexity of K-mean is $O(nkd_i)$, $n$ is the number of data points, $d$ is the number of features, $k$ is the number of clusters, $i$ is iterations.

Agglomerative clustering : 1). Maintain a set of clusters, initially, each instance in its own cluster.

produce not one clustering, but a dendrogram tree.

2). Repeat : i) Pick the two closest clusters, merge them into a new cluster.
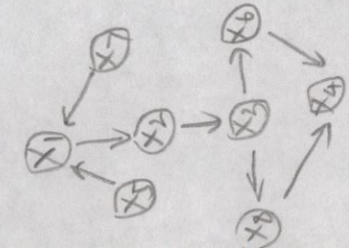              ii) stop when there is only one cluster left.

VC dimension of hyperplane in d-dimensional space is d+1.

The VC dimension of all axis parallel rectangles is $\geq 4$.

no 5 instances can be shattered since there can be at most $2^4 = 16$.

A distinct extreme points and these 4 points cannot be included without including any possible 5th point.

no point
1 point
2 point
3 point
4 point

$VC(H) = 2n$, $n$ is dimension of axis-parallel hyper-rectangle.

Hypothesis space $|H| = n^m$, $n$ is the possible range of variable, $m$ is the number of variables.

Elimination : $x_1, x_2, x_3, x_5, x_6, x_7, x_8$, to compute $P(x_4 = x_4)$.

$x_1$, $\phi(x_7, x_5, x_2) = \sum_{x_1} P(x_1|x_7, x_5) \cdot P(x_2|x_1)$    exp(4)

$x_2$, $\phi(x_5, x_7, x_3) = \sum_{x_2} \phi(x_7, x_5, x_2) \cdot P(x_3|x_2)$    exp(4)

$x_3$, $\phi(x_5, x_7, x_6, x_8) = \sum_{x_3} \phi(x_5, x_7, x_3) \cdot P(x_6|x_3) \cdot P(x_8|x_3)$    exp(5)

$x_5$, $\phi(x_7, x_6, x_8) = \sum_{x_5} \phi(x_5, x_7, x_6, x_8) \cdot P(x_5)$    exp(4).

$x_6$, $\phi(x_7, x_8)$
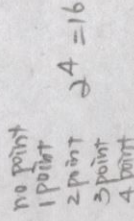                            $= \sum_{x_6} \phi(x_7, x_6, x_8) \cdot P(x_4 = x_4 | x_6, x_8)$    exp(3)

$x_7$, $\phi(x_8) = \sum_{x_7} \phi(x_7, x_8) P(x_7)$    exp(2)

$x_8$, $P(x_4 = x_4) = \sum_{x_8} \phi(x_8)$    exp(1).

Time complexity = $O(7 \cdot exp(5))$ , Space complexity = $O(7 \cdot exp(4))$.
                                                          size of the largest function created.

$\Phi(X_8) = \sum_{X_7} \cdot \Phi(X_7, X_8) P(X_7)$     $\exp(2)$

$X_8, \quad P(X_4 = X_4) = \sum_{X_8} \Phi(X_8)$     $\exp(1)$.

Time complexity = $O(7 \cdot \exp(5))$, Space complexity: $O(7 \cdot \exp(4))$.
Size of the largest function created.

BN 1 : have no-edge, BN2 : have an edge.
$X_1$ and $X_2$ independent    $X_1$ and $X_2$ not independent.

Assume no noisy, then 1-NN is 100% true because the class of each example is the class of its nearest neighbor, that is itself.

Some functions, e.g. the parity function requires exponential number of hidden nodes, if only one hidden layer is used.

Each node in the decision tree divides the examples into two non-empty sets, n training examples, then the number of nodes from the leaf to the root is bounded by "n".