

Instance Based Learning

- given: an instance $x^{(0)}$ to classify
 - find the k training-set instances $(x^{(1)}, y^{(1)}) \dots (x^{(k)}, y^{(k)})$ that are most similar to $x^{(0)}$
 - return the class value
 $\hat{y} \leftarrow \arg \max_{i \in \text{values}(Y)} \sum_{j=1}^k \delta(x^{(j)}, y^{(j)})$
 $\delta(a, b) = \begin{cases} 1 & \text{if } a = b \\ 0 & \text{otherwise} \end{cases}$
- incremental deletion**
start with all training instances in memory
for each training instance $(x^{(0)}, y^{(0)})$
 if other training instances provide correct classification for $(x^{(0)}, y^{(0)})$
 delete it from the memory
- incremental growth**
start with an empty memory
for each training instance $(x^{(0)}, y^{(0)})$
 if other training instances in memory don't correctly classify $(x^{(0)}, y^{(0)})$
 add it to the memory

Find nearest neighbor in k-d tree

```
NearestNeighbor(instance x(0))
PO = () // minimizing priority queue
best_dist = ∞ // smallest distance seen so far
PO.push(root, 0)
while PO is not empty
  (node, bound) = PO.pop();
  if (bound ≥ best_dist)
    return best_node.instance // nearest neighbor found
  dist = distance(x(0), node.instance)
  if (dist < best_dist)
    best_dist = dist
    best_node = node
  if (q[node feature] - node.threshold > 0)
    PQ.push(node.left, x(0)[node.feature] - node.threshold)
    PQ.push(node.right, 0)
  else
    PQ.push(node.left, 0)
    PQ.push(node.right, node.threshold - x(0)[node.feature])
return best_node.instance
```

Decision Trees

```
MakeSubtree(set of training instances D)
C = DetermineCandidateSplits(D)
if stopping criteria met
  make a leaf node N
  determine class label/probabilities for N
else
  make an internal node N
  S = FindBestSplit(D, C)
  for each outcome k of S
    Dk = subset of instances that have outcome k
    kth child of N = MakeSubtree(Dk)
return subtree rooted at N
```

$$H(Y) = - \sum_y P(y) \log_2 P(y)$$

$$H(Y|X) = \sum_x P(X=x) H(Y|X=x)$$

$$H(Y|X=x) = - \sum_y P(Y=y|X=x) \log_2 P(Y=y|X=x)$$

$$\sum_{i=1}^{|D|} (y_i - \hat{y}_i)^2 = \sum_{L \in \text{leaves}} \sum_{i \in L} (y_i - \hat{y}_i)^2$$

Pruning in C4.5

- split given data into training and tuning (validation) sets
- grow a complete tree
- do until further pruning is harmful
 - evaluate impact on tuning-set accuracy of pruning each node
 - greedily remove the one that most improves tuning-set accuracy

Methodology

$$\text{TPR (recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

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

$$\text{error}_D(h) \pm z_\alpha \sqrt{\frac{\text{error}_D(h)(1 - \text{error}_D(h))}{n}}$$

$$\bar{\delta} = \frac{t}{\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (\delta_i - \bar{\delta})^2}}$$

Learning curve

- given training/test set partition
- for each sample size s on learning curve
 - (optionally) repeat n times
 - randomly select s instances from training set
 - learn model
 - evaluate model on test set to determine accuracy a
 - plot (s, a) or $(s, \text{avg. accuracy and error bars})$

ROC curve

```
let  $\{(y^{(1)}, c^{(1)}) \dots (y^{(n)}, c^{(n)})\}$  be the test-set instances sorted according to predicted confidence  $c^{(i)}$  that each instance is positive
let num_pos, num_neg be the number of negative/positive instances in the test set
TP = 0, FP = 0
last_TP = 0
for i = 1 to n
  // find thresholds where there is a pos instance on high side, neg instance on low side
  if (i > 1) and ( $c^{(i)} \neq c^{(i-1)}$ ) and ( $y^{(i)} == \text{neg}$ ) and ( $TP > \text{last\_TP}$ )
    FPR = FP / num_neg, TPR = TP / num_pos
    output (FPR, TPR) coordinate
    last_TP = TP
  if  $y^{(i)} == \text{pos}$ 
    + +TP
  else
    + +FP
FPR = FP / num_neg, TPR = TP / num_pos
output (FPR, TPR) coordinate
```

Bayes Network

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

$$P(b|j, m) = \frac{P(b, j, m)}{P(j, m)} = \frac{P(b, j, m)}{P(b, j, m) + P(\neg b, j, m)}$$

$$\text{InfoGain}(D, S) = H_D(Y) - H_D(Y|S)$$

$$\text{SplitInfo}(D, S) = - \sum_{k \in \text{outcomes}(S)} \frac{|D_k|}{|D|} \log_2 \left(\frac{|D_k|}{|D|} \right)$$

$$\text{GainRatio}(D, S) = \frac{\text{InfoGain}(D, S)}{\text{SplitInfo}(D, S)}$$

$$\text{MAP (Laplace): } P(X=x) = \frac{n_x + 1}{\sum_{v \in \text{values}(X)} (n_v + 1)}$$

$$\text{m-estimates: } P(X=x) = \frac{n_x + p_x m}{(\sum_{v \in \text{values}(X)} n_v) + m}$$

$$\text{M-step: } P(a|b, e) = \frac{E\#(a \wedge b \wedge e)}{E\#(b \wedge e)}$$

$$I(X, Y) = \sum_x \sum_y P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$$

$$P(X_1, \dots, X_n, Y) = P(Y) \prod_{i=1}^n P(X_i | Y)$$

$$P(Y = y | X) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{\sum_{y'} P(y') \prod_{i=1}^n P(x_i | y')}$$

$$I(X_i, X_j | Y) = \sum_{x_i} \sum_{x_j} \sum_y P(x_i, x_j, y) \log_2 \frac{P(x_i, x_j | y)}{P(x_i | y)P(x_j | y)}$$

$$P(Y = y | \mathbf{x}) = \frac{P(y)P(\mathbf{x}|y)}{\sum_{y'} P(y')P(\mathbf{x}|y')}$$

TAN algorithm

- compute weight $I(X_i, X_j | Y)$ for each possible edge (X_i, X_j) between features
- find maximum weight spanning tree (MST) for graph over $X_1 \dots X_n$
- assign edge directions in MST
- construct a TAN model by adding node for Y and an edge from Y to each X_i

Markov Network

$$P(\mathbf{V}) = \frac{1}{Z} \prod_k \phi_k(\mathbf{v})$$

$$Z = \sum_{\mathbf{v} \in \mathbf{V}} \prod_k \phi_k(\mathbf{v})$$

Neural Network

$$E(\mathbf{w}) = \frac{1}{2} \sum_{d \in D} (y^{(d)} - o^{(d)})^2$$

$$E(\mathbf{w}) = \sum_{d \in D} -y^{(d)} \ln(o^{(d)}) - (1 - y^{(d)}) \ln(1 - o^{(d)})$$

$$E(\mathbf{w}) = - \sum_{d \in D} \sum_{i=1}^{\#\text{class}} y_i^{(d)} \ln(o_i^{(d)})$$

$$\text{net} = w_0 + \sum_i w_i x_i$$

$$f(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad (\text{Sigmoid})$$

$$f(\text{net}) = \tanh(x) = \frac{2}{1 + e^{-2\text{net}}} - 1 \quad (\text{Tanh})$$

$$f(\text{net}) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases} \quad (\text{ReLU})$$

$$f(\text{net}) = \frac{e^{\text{net}_j}}{\sum_j e^{\text{net}_j}} \quad (\text{softmax})$$

$$\Delta \mathbf{w} = -\eta \nabla_{\mathbf{w}} E(\mathbf{w})$$

$$\Delta w_{ji} = -\eta \delta_j o_i$$

$$\delta_j = -\frac{\partial E}{\partial \text{net}_j} = \text{e.g. } o_j(1 - o_j) \sum_k \delta_k w_{kj}$$

given: network structure and a training set $D = \{(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})\}$

initialize all weights in \mathbf{w} to small random numbers

until stopping criteria met do

 initialize the error $E(\mathbf{w}) = 0$

 for each $(x^{(d)}, y^{(d)})$ in the training set

 input $x^{(d)}$ to the network and compute output $o^{(d)}$

 increment the error $E(\mathbf{w}) = E(\mathbf{w}) + \frac{1}{2} (y^{(d)} - o^{(d)})^2$

 calculate the gradient

$$\nabla E(\mathbf{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

 update the weights

$$\Delta \mathbf{w} = -\eta \nabla E(\mathbf{w})$$

Step 1: $w^{(t+1)} = \mu v^{(t)} - \eta \cdot \frac{\partial E}{\partial w}|_{w^{(t)}}$

Interpret v as velocity, $1 - \mu$ as friction, v stores the moving average of the gradients.

Step 2: $w^{(t+1)} = w^{(t)} + v^{(t+1)}$

Step 1: $v^{(t+1)} = \mu v^{(t)} - \eta \cdot \frac{\partial E}{\partial w}|_{w^{(t)} + \mu v^{(t)}}$

Step 2: $w^{(t+1)} = w^{(t)} + v^{(t+1)}$

Step 1: $G^{(t+1)} = G^{(t)} + \left(\frac{\partial E}{\partial w}\right)_{w^{(t)}}^2$
 G is the sum of the square gradient

Step 2: $w^{(t+1)} = w^{(t)} - \frac{\eta}{\sqrt{G^{(t+1)} + 10^{-7}}} \odot \frac{\partial E}{\partial w}|_{w^{(t)}}$
Scale the learning rate by $1/\sqrt{G}$ and follow the negative gradient direction

Step 1: $G^{(t+1)} = \gamma G^{(t)} + (1-\gamma) \left(\frac{\partial E}{\partial w}\right)_{w^{(t)}}^2$
 G is a moving exponential average of the square gradient.

Step 2: $w^{(t+1)} = w^{(t)} - \frac{\eta}{\sqrt{G^{(t+1)} + 10^{-7}}} \odot \frac{\partial E}{\partial w}|_{w^{(t)}}$
Scale the learning rate by $1/\sqrt{G}$ and follow the negative gradient direction

Step 1: $m^{(t+1)} = \beta_1 m^{(t)} + (1-\beta_1) \frac{\partial E}{\partial w}|_{w^{(t)}}$
 m is the exponential moving average of the gradients (similar to speed in the momentum method)

Step 2: $v^{(t+1)} = \beta_2 v^{(t)} + (1-\beta_2) \cdot \left(\frac{\partial E}{\partial w}\right)_{w^{(t)}}^2$
 v is the exponential moving average of the square gradients

Step 3: $w^{(t+1)} = w^{(t)} - \frac{\eta}{\sqrt{v^{(t+1)} + 10^{-9}}} \odot m^{(t+1)}$

$$E(\theta^{(D)}) = -\frac{1}{2} E_{x \sim p_{data}} \log D(\mathbf{x}) - \frac{1}{2} E_z \log (1 - D(G(\mathbf{z})))$$

$$E(\theta^{(G)}) = -E(\theta^{(D)})$$

$$\min_G \max_D E_{x \sim p_{data}} \log D(\mathbf{x}) + E_z \log (1 - D(G(\mathbf{z})))$$

Linear Models

Ridge regression

$$\hat{\beta} = (X^T X - \lambda I)^{-1} X^T Y$$

LASSO coordinate descent

repeat until convergence

$$\begin{aligned} \text{for } j = 1 \dots n \\ \rho_j &= \sum_i x_j^{(i)} \left(y^{(i)} - \sum_{k \neq j} \beta_k x_k^{(i)} \right) \\ z_j &= \sum_i (x_j^{(i)})^2 \\ \beta_j &= \frac{1}{z_j} S(\rho_j, \lambda) \end{aligned}$$

$$S(\rho, \lambda) = \begin{cases} \rho - \lambda & \text{if } \rho > \lambda \\ 0 & \text{if } -\lambda \leq \rho \leq \lambda \\ \rho + \lambda & \text{if } \rho < -\lambda \end{cases}$$

Ensembles

Bagging

learning:
given: learner L , training set $D = \{ \langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle \}$
for $i \leftarrow 1$ to T do
 $D_i \leftarrow m$ instances randomly drawn with replacement from D
 $h_i \leftarrow$ model learned using L on D_i

AdaBoost

given: learner L , # stages T , training set $D = \{ \langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle \}$

for all i : $w_i(i) \leftarrow 1/m$ // initialize instance weights
for $t \leftarrow 1$ to T do
for all i : $p_t(i) \leftarrow w_t(i) / (\sum_j w_t(j))$ // normalize weights
 $h_t \leftarrow$ model learned using L on D and p_t
 $e_t \leftarrow \sum_i p_t(i)(1 - h_t(x^{(i)}, y^{(i)}))$ // calculate weighted error
if $e_t > 0.5$ then
 $T \leftarrow t - 1$
break
 $\beta_t \leftarrow e_t / (1 - e_t)$
for all i where $h_t(x^{(i)}) = y^{(i)}$ // down-weight correct examples
 $w_{t+1}(i) \leftarrow w_t(i) \beta_t$

return:
 $h(x) = \arg \max_y \sum_{i=1}^T \left(\log \frac{1}{\beta_i} \right) \delta(h_i(x), y)$

Random Forest

given: candidate feature splits F ,
training set $D = \{ \langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle \}$
for $i \leftarrow 1$ to T do
 $D_i \leftarrow m$ instances randomly drawn with replacement from D
 $h_i \leftarrow$ randomized decision tree learned with F, D_i

randomized decision tree learning:

to select a split at a node:
 $R \leftarrow$ randomly select (without replacement) f feature splits from F
(where $f << |F|$)
choose the best feature split in R
do not prune trees

Support Vector Machines

$$h(x) = \begin{cases} +1 & \text{if } \mathbf{w}^T \mathbf{x} + b > 0 \\ -1 & \text{o/w} \end{cases}$$

$$\text{margin}_D(h) = \frac{1}{2} \mathbf{w}^T (\mathbf{x}_+ - \mathbf{x}_-) = \frac{1}{\|\mathbf{w}\|_2}$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad (\text{hard-margin})$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \ \forall i$$

$$\min_{\mathbf{w}, b, \xi^{(i)}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi^{(i)} \quad (\text{soft-margin})$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi^{(i)}, \xi^{(i)} \geq 0 \ \forall i$$

$$\max_{\alpha_1, \dots, \alpha_m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{j=1}^m \sum_{i=1}^m \alpha_j \alpha_k y^{(j)} y^{(k)} (\mathbf{x}^{(j)} \cdot \mathbf{x}^{(k)})$$

$$\text{s.t. } \alpha_i \geq 0, \sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (\text{h-m dual})$$

$$k(x, z) = (x \cdot z)^d \quad (\text{polynomial d})$$

$$k(x, z) = (x \cdot z + 1)^d \quad (\text{polynomial up-to d})$$

$$k(x, z) = \exp(-\gamma \|x - z\|^2) \quad (\text{RBF})$$

SMO algorithm

• Input: C , kernel, kernel parameters, ϵ

• Initialize b and all α 's to 0

• Repeat until KKT satisfied (to within ϵ):

– Find an example **e1** that violates KKT (prefer unbound examples here, choose randomly among those)

– Choose a second example **e2**. Prefer one to maximize step size (in practice, faster to just maximize $|E_1 - E_2|$). If that fails to result in change, randomly choose unbound example. If that fails, randomly choose example. If that fails, re-choose **e1**.

– Update α_1 and α_2 in one step, compute new threshold b

$$\alpha_i = 0 \iff y_i u_i \geq 1$$

$$0 < \alpha_i < C \iff y_i u_i = 1$$

$$\alpha_i = C \iff y_i u_i \leq 1$$

• Given examples **e1** and **e2**, set

$$\alpha_2^{\text{new}} = \alpha_2 + \frac{y_2(E_1 - E_2)}{\eta}$$

where $\eta = K(x_1, x_1) + K(x_2, x_2) - 2K(x_1, x_2)$

• Clip this value in the natural way:

$$\alpha_2^{\text{new, clipped}} = \begin{cases} H & \text{if } \alpha_2^{\text{new}} \geq H, \\ \alpha_2^{\text{new}} & \text{if } \alpha_2^{\text{new}} < H; \\ L & \text{if } \alpha_2^{\text{new}} \leq L. \end{cases}$$

if $y_1 = y_2$ then:

$$\begin{cases} L = \max(0, \alpha_2 + \alpha_1 - C) \\ H = \min(C, \alpha_2 + \alpha_1) \end{cases}$$

otherwise:

$$\begin{cases} L = \max(0, \alpha_2 - \alpha_1) \\ H = \min(C + \alpha_2 - \alpha_1, C) \end{cases}$$

• Set $\alpha_1^{\text{new}} = \alpha_1 + s(\alpha_2 - \alpha_2^{\text{new, clipped}})$ where $s = y_1 y_2$

$$E[(y - f(x; D))^2] = E[(y - E[y|x])^2 | x, D]$$

$$+ (f(x; D) - E[y|x])^2$$

$h \in H$ overfits if $\exists h' \in H$ s.t. $\text{error}(h) > \text{error}(h')$
and $\text{error}_D(h) < \text{error}_D(h')$

Reinforcement Learning

$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t E[r_t]$$

$$\pi^* = \arg \max_{\pi} V^\pi(s) \quad \forall s$$

$$Q(s, a) = E[r(s, a)] + \gamma E_{s'|s,a}[V^*(s')]$$

$$\pi^*(s) = \arg \max_a Q(s, a)$$

$$P(a_i|s) = \frac{c^{\hat{Q}(s, a_i)}}{\sum_j c^{\hat{Q}(s, a_j)}}$$

Q Learning (Deterministic)

for each s, a initialize table entry $\hat{Q}(s, a) \leftarrow 0$

observe current state s

do forever

select an action a and execute it

receive immediate reward r

observe the new state s'

update table entry

$$\hat{Q}(s, a) \leftarrow r + \gamma \max_a \hat{Q}(s', a)$$

$$s \leftarrow s'$$

Q Learning (Non-deterministic)

for each s, a initialize table entry $\hat{Q}(s, a) \leftarrow 0$

observe current state s

do forever

select an action a and execute it

receive immediate reward r

observe the new state s'

update table entry

$$\hat{Q}(s, a) \leftarrow (1 - \alpha_s) \hat{Q}_{s-1}(s, a) + \alpha_s [r + \gamma \max_a \hat{Q}_{s-1}(s', a)]$$

$$s \leftarrow s'$$

AlphaZero

Initialize DNN f_θ

Repeat Forever

Play Game

Update θ

Play Game:

Repeat Until Win or Lose:

From current state S , perform MCTS

Estimate move probabilities π by MCTS

Record (S, π) as an example

Randomly draw next move from π

Bias, Variance, & Overfitting

$$E_D[(f(x; D) - E[y|x])^2] = (E_d[f(x; D)] - E[y|x])^2$$

$$+ E_D[(f(x; D) - E_d[f(x; D)])^2]$$

Update θ :

Let z be previous game outcome (+1 or -1)

Sample from last game's examples (S, π, z)

Train DNN f_θ on sample to get new θ