

- Model Selection
 - Cross Validation
 - Split the data set \mathcal{D} into K subsets $\mathcal{D}^{(i)}$ (called folds)
 - For $i = 1, \dots, K$, train $h^{(i)}$ using all data but the i -th fold ($\mathcal{D} \setminus \mathcal{D}^{(i)}$).
 - Cross-validation error by averaging all validation errors $\hat{\epsilon}_{D(w)}(h^{(i)})$

- KNN
 - 误差界估计
 - Normalization
 - Mean $\mu_j = \frac{1}{N} \sum_{i=1}^N x_{ij}$
 - Variance $\sigma_j = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_{ij} - \mu_j)^2}$
 - Normalize the feature into a new one:

$$\hat{x}_{ij} \leftarrow \frac{x_{ij} - \mu_j}{\sigma_j}$$

- Distance
- Cosine Distance: $\cos(x_i, x_j) = \frac{\langle x_i, x_j \rangle}{\|x_i\| \|x_j\|}$
 - Minkowski Distance: $D_p(x_i, x_j) = \sqrt[p]{\sum_{k=1}^d |x_{ik} - x_{jk}|^p}$
 - Mahalanobis Distance: $D_M(x_i, x_j) = \sqrt{(x_i - x_j)^T M (x_i - x_j)}$

- Weighted KNN
- Distance-weighted Classification:

$$\hat{y} = \operatorname{argmax}_{c \in Y} \sum_{x' \in \operatorname{KNN}_K(x)} \frac{1}{D(x, x')^2}$$
 - Distance-weighted Regression:

$$\hat{y} = \sum_{x' \in \operatorname{KNN}_K(x)} \frac{1}{D(x, x')^2} y'$$

- Nearest Centroid Classifier
- Given dataset $\{x_i, y_i\}_{i=1}^N$ with class label $y_i \in Y$:
 - Compute per-class centroids $\mu_c = \frac{1}{|N_c|} \sum_{y_i=c} x_i$.
 - Prediction function: $\hat{y} = \operatorname{argmin}_{c \in Y} \|x - \mu_c\|$
 - $O(Nd)$ training and $O(Cd)$ testing.
 - Faster!
 - Gaussian assumption for each class.
 - Cannot tackle complex distributions.
 - Very good for Few-Shot Learning (FSL).

- Complexity: n examples, d dimensions
- For naïve KNN:
 - $O(1)$ for training.
 - $O(nd)$ to find k closest examples.
 - For KNN with k -d Tree:
 - $O(dn \log n)$ for training (build k -d tree)
 - $O(2^d \log n)$ on average when query.
 - Use k -d tree only when $n \gg 2^d$.

- Summary
- When to use:
 - Few attributes per instance (expensive computation)
 - Lots of training data (curse of dimensionality)
 - Advantages:
 - Agnostically learn complex target functions
 - Do not lose information (store original data)
 - Data number can be very large (big prof)
 - Class number can be very large (biggest prof)
 - All other ML algorithms may fail here!
 - Disadvantages:
 - Slow at inference time (acceleration a must)
 - Ineffective in high dimensions (curse of dimensionality)
 - Fooled easily by irrelevant attributes (feature engineering crucial)

- Linear Regression
 - Analytic Sol
 - Computing the gradient of $\hat{\epsilon}(w)$ w.r.t. w and setting it to zero yields the optimal parameter w^* :

$$\begin{aligned} \hat{\epsilon}(w) &= \|Xw - y\|^2 \\ \nabla_w \hat{\epsilon}(w) &= 2X^T(Xw - y) = 0 \\ &\Rightarrow X^T Xw = X^T y \\ &\Rightarrow w = (X^T X)^{-1} X^T y \end{aligned}$$

- Optimization: SGD
- In each iteration $t (\leq T)$:
 - Randomly sample a minibatch of $m \ll n$ points $\{(x_i, y_i)\}_{i=1}^m$
 - Set $J^t(w^t) = \frac{1}{m} \sum_{i=1}^m \ell(w^t; x_i, y_i)$
 - Compute gradient on minibatch:

$$\Delta^t = \nabla_w J^t(w^t)$$
 - Update parameters with learning rate η :

$$w^{t+1} = w^t - \eta \Delta^t$$

- For linear regression: $w^{t+1} \leftarrow w^t - 2\eta X_m^T (X_m w^t - y)$
 - Computational complexity $O(dmT)$. What if large dimension d ?
 - Theoretically, convergence rate is $O(1/\sqrt{T})$ under convex condition.
- Regularization
- Norm-regularization will make the hypothesis smooth at any point.
- Linear Classification
 - Logistic Regression
 - To map the output of $h(x)$ into $[0, 1]$, we use sigmoid function:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

- Sigmoid maps \mathbb{R} to $[0, 1]$.
 - $t \rightarrow +\infty, \sigma(t) \rightarrow 1$.
 - Set $\sigma(h(x)) = p(y = 1|x)$ as the probability to label x as $y = 1$.
- $$\hat{\epsilon}(w) = - \sum_{i=1}^n [y_i \log \sigma(h_w(x)) + (1 - y_i) \log [1 - \sigma(h_w(x))]] + \lambda \sum_{j=1}^d w_j^2$$

- SoftMax Regression
- Softmax function normalizes multiple outputs in a probability vector:

$$p(y = c|x; W) = \frac{\exp(w_c^T x)}{\sum_{r=1}^C \exp(w_r^T x)}$$
 - Loss for each data point (x_i, y_i) :

$$\ell(h(x_i), y_i) = \begin{cases} -\log \frac{\exp(w_1^T x_i)}{\sum_{r=1}^C \exp(w_r^T x_i)}, & y_i = 1 \\ -\log \frac{\exp(w_2^T x_i)}{\sum_{r=1}^C \exp(w_r^T x_i)}, & y_i = 2 \\ \vdots \\ -\log \frac{\exp(w_C^T x_i)}{\sum_{r=1}^C \exp(w_r^T x_i)}, & y_i = C \end{cases}$$

- SVM
 - SM SVM
 - Original form

$$\begin{aligned} \min_{w,b,\xi} & \frac{1}{2} \|w\|^2 \\ \text{s.t. } & y_i(w \cdot x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, \sum_{i=1}^n \xi_i \leq n', 1 \leq i \leq n \end{aligned}$$
 - Lagrangian equivalent form

$$\begin{aligned} \min_{w,b,\xi} & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & y_i(w \cdot x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, 1 \leq i \leq n \end{aligned}$$

- Dual Problem

- Primal Problem

$$\begin{aligned} \min_{w,b} & \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \max[0, 1 - y_i(w \cdot x_i + b)] \\ \text{s.t. } & g_j(x) \leq 0 \text{ for } j = 1, \dots, J \\ & h_k(x) = 0 \text{ for } k = 1, \dots, K \end{aligned}$$
 - With Lagrange multipliers λ, μ , the Lagrangian function is defined as

$$L(x, \lambda, \mu) = f(x) + \sum_{j=1}^J \lambda_j g_j(x) + \sum_{k=1}^K \mu_k h_k(x)$$

- KKT conditions
- KKT conditions: for $1 \leq j \leq J$
 - Primal feasibility: $g_j(x) \leq 0, h_k(x) = 0$
 - Dual feasibility: $\lambda_j \geq 0$
 - Complementary slackness: $\lambda_j g_j(x) = 0$
 - The gradient of the Lagrangian function w.r.t. x vanishes to 0:

$$\nabla_x f(x) + \sum_{j=1}^J \lambda_j \nabla_x g_j(x) + \sum_{k=1}^K \mu_k \nabla_x h_k(x) = 0$$
 - If x^*, λ^*, μ^* satisfy KKT for a convex problem, then x^* is an optimum:

$$f(x^*) = L(x^*, \lambda^*, \mu^*)$$
 - x^* is optimal if and only if there exist λ, μ satisfying KKT conditions.

- Dual Problem
 - Define the primal objective as

$$\Pi(x) = \max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} L(x, \lambda, \mu)$$
 - Define the dual objective as

$$\Gamma(\lambda, \mu) \triangleq \min_{x \in \mathbb{R}^d} L(x, \lambda, \mu)$$
 - Weak duality:

$$\max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} \Gamma(\lambda, \mu) \leq \min_{x \in \mathbb{R}^d} \Pi(x)$$
 - Dual Problem:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} & \Gamma(\lambda, \mu) \\ \text{s.t. } & \lambda_j \geq 0 \text{ for } j = 1, \dots, J \end{aligned}$$

- Soft-Margin SVM Dual Problem
 - Lagrangian function (with $2n$ inequality constraints):

$$L(w, b, \alpha, \xi, \mu)$$
 - $$\begin{aligned} &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(w \cdot x_i + b)) - \sum_{i=1}^n \mu_i \xi_i \\ &\alpha_i \geq 0, \mu_i \geq 0, i = 1, \dots, n \end{aligned}$$
 - After solving for α , we can solve for $w = \sum_{i=1}^n \alpha_i y_i x_i$.

- Primal Problem: SGD
- Theorem: Soft-SVM is equivalent to Regularized Risk Minimization:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + \frac{C}{n} \sum_{i=1}^n \max[0, 1 - y_i(w \cdot x_i + b)]$$

- SGD for Soft-SVM:
 - In each iteration $t (\leq T)$:
 - Choose i uniformly at random from $1, \dots, n$ (not in mini-batches)
 - The subgradient $g_i^t = \begin{cases} w^t & \text{if } y_i(w^t \cdot x_i + b) \geq 1 \\ w^t - C y_i x_i & \text{otherwise} \end{cases}$
 - Update $w^{t+1} \leftarrow w^t - \eta g_i^t$
 - Output $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^t$

- Kernel SVM
- Commonly used kernel functions for vector data:

Linear	$k(x_1, x_2) = x_1 \cdot x_2$
Polynomial	$k(x_1, x_2) = (x_1 \cdot x_2)^d$
Polynomial (with constant terms)	$k(x_1, x_2) = (x_1 \cdot x_2 + 1)^d$
RBF (Gaussian kernel)	$k(x_1, x_2) = \exp\left(-\frac{\ x_1 - x_2\ ^2}{2\sigma^2}\right)$
Hyperbolic Tangent	$k(x_1, x_2) = \tanh(\beta x_1 \cdot x_2 + c)$

- From kernel functions k_1, k_2 , we can construct new kernel functions:
 - $k'(x_1, x_2) = k_1 \otimes k_2(x_1, x_2) = k_1(x_1, x_2) k_2(x_1, x_2)$
 - For any function $g: \mathcal{X} \rightarrow \mathbb{R}$, $k'(x_1, x_2) = g(x_1) k_1(x_1, x_2) g(x_2)$
- Theorem (Mercer): If $k(\cdot, \cdot)$ is a symmetric function on space $\mathcal{X} \times \mathcal{X}$, then k is a kernel function \Leftrightarrow For any input set (x_1, x_2, \dots, x_m) ,

$$K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_m) \\ \vdots & \ddots & \vdots \\ k(x_m, x_1) & \dots & k(x_m, x_m) \end{pmatrix}$$
 The kernel matrix is semi-definite ($K \succcurlyeq 0$, i.e. $x^T K x \geq 0, \forall x$).
- For kernel functions k_1, k_2, \dots, k_s , and $\gamma_1, \gamma_2, \dots, \gamma_s \geq 0$:
 - $\sum_{s=1}^s \gamma_s k_s$ is also (multi-)kernel function, because $\sum_{s=1}^s \gamma_s K_s \succcurlyeq 0$.
- Dual Problem of Soft-SVM:

$$\begin{aligned} \max_{\alpha} & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{s.t. } & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, 1 \leq i \leq n. \end{aligned}$$

- After solving for α , we can solve for $w = \sum_{i=1}^n \alpha_i y_i \Phi(x_i)$.
- Testing: $f(x) = w \cdot \Phi(x) + b = \sum_{i=1}^n \alpha_i y_i k(x_i, x) + b$.
- During training and testing, only need to compute $k(x_i, x_j)$.

- Learning Theory
 - Bias-Variance Decomposition
 - Theorem: $\mathcal{E}_{L_2}(x) = \text{Var}(y|x) + \text{Bias}[h_D(x)|x]^2 + \text{Var}_D[h_D(x)|x]$

- Approximation Error and Estimation Error
- Diagram showing Error vs. Complexity. The total error is the sum of Bias² (decreasing), Variance (increasing), and Irreducible error (constant). The Bias-Variance decomposition is shown as $\text{Bias}^2[h_D(x)|x]^2 + \text{Var}_D[h_D(x)|x]$.

- PAC Learning
- General equality: given a target function f , for any $h \in \mathcal{H}$,

$$\mathcal{E}(h) - \mathcal{E}^*(f) = [\mathcal{E}(h) - \mathcal{E}(h^*)] + [\mathcal{E}(h^*) - \mathcal{E}^*(f)]$$
 estimation approximation

• A hypothesis space \mathcal{H} is **PAC-learnable** if there exists an algorithm \mathcal{A} and a polynomial **poly()**, such that for any $\epsilon > 0, \delta > 0$, for all distributions D on \mathcal{X} and for any target hypothesis $h \in \mathcal{H}$, the following holds for **sample complexity** $n \geq \text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, |\mathcal{H}|)$:

$$P_{\mathcal{D}_n \sim D^n} \left[\mathcal{E}(h_{\mathcal{D}_n}) - \min_{h \in \mathcal{H}} \mathcal{E}(h) \geq \epsilon \right] \leq \delta$$

$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \mathcal{E}(h)$

Approximately correct

Probably correct

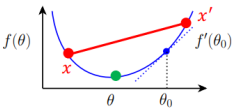


d) Probability Tools

Basic

- **Union bound:** $P(A \vee B) \leq P(A) + P(B)$.
- **Inversion:** if $P(X \geq \epsilon) \leq f(\epsilon)$, then for any $\delta > 0$, with probability at least $1 - \delta$, $X \leq f^{-1}(\delta)$.
 - Useful to interchange ϵ and δ : $P(|\hat{\mathcal{E}}_{\mathcal{D}_n}(h) - \mathcal{E}(h)| \geq \epsilon) \leq \delta$.
- **Jensen's inequality:** if f is convex, then $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$.
 - Convex function: f is convex if $\forall \lambda \in [0,1], x, x' \in \text{dom}(f)$:

$$f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x')$$
 - $f(x') \geq f(x) + g^T(x' - x)$
 - g : Subgradient.



Concentration Inequalities

Markov's Inequality

• **Theorem:** If Z is a nonnegative random variable, then for any $\epsilon > 0$:

$$P(Z \geq \epsilon) \leq \frac{\mathbb{E}[Z]}{\epsilon}$$

Proof:

Chebyshev's Inequality

$$P(V \geq \epsilon) \leq P(\phi(V) \geq \phi(\epsilon)) \leq \frac{\mathbb{E}[\phi(V)]}{\phi(\epsilon)}$$

Chernoff Bound

• **Theorem:** If we take $\phi(t) = e^{\lambda t}$, we have for any $\lambda > 0$:

$$P(V \geq \epsilon) \leq \frac{\mathbb{E}[e^{\lambda V}]}{e^{\lambda \epsilon}}$$

• If we take $V = Z - \mathbb{E}[Z]$, we have:

$$P(Z - \mathbb{E}[Z] \geq \epsilon) \leq \frac{\mathbb{E}[e^{\lambda(Z - \mathbb{E}[Z])}]}{e^{\lambda \epsilon}}$$

• If we set Z as the **sum** of **i.i.d.** random variables: $Z = \sum_{i=1}^n X_i$

$$P\left(\sum_{i=1}^n (X_i - \mathbb{E}X_i) \geq \epsilon\right) \leq \frac{\mathbb{E}[e^{\lambda \sum_{i=1}^n (X_i - \mathbb{E}X_i)}]}{e^{\lambda \epsilon}} = \frac{\prod_{i=1}^n \mathbb{E}[e^{\lambda (X_i - \mathbb{E}X_i)}]}{e^{\lambda \epsilon}}$$

Hoeffding's Lemma

• **Lemma:** Let V be a **bounded** random variable, $\mathbb{E}[V] = 0$ and $a \leq V \leq b$ with $b > a$. Then for any $\lambda > 0$:

$$\mathbb{E}[e^{\lambda V}] \leq e^{\frac{\lambda^2 (b-a)^2}{8}}$$

• Plug in the optimal $\lambda = 4\epsilon / \sum_{i=1}^n (b_i - a_i)^2 \Rightarrow$ Hoeffding's inequality:

$$P\left(\sum_{i=1}^n (X_i - \mathbb{E}X_i) \geq \epsilon\right) \leq \frac{\prod_{i=1}^n \mathbb{E}[e^{\lambda (X_i - \mathbb{E}X_i)}]}{e^{\lambda \epsilon}} \leq \frac{e^{\sum_{i=1}^n \frac{\lambda^2 (b_i - a_i)^2}{8}}}{e^{\lambda \epsilon}} \stackrel{\text{min}}{\implies} \exp\left[-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right]$$

Chernoff Bound

Hoeffding's Lemma

Optimal λ

• If we replace X_i with $-X_i \in [-b_i, -a_i]$, we have

$$P\left(\sum_{i=1}^n (X_i - \mathbb{E}X_i) \leq -\epsilon\right) \leq \exp\left[-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right]$$

• **Theorem** (Hoeffding's Inequality): for $a_i \leq X_i \leq b_i$, **finite-sample bound**

$$P\left(\left|\sum_{i=1}^n (X_i - \mathbb{E}X_i)\right| \geq \epsilon\right) \leq 2 \exp\left[-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right]$$

McDiarmid's Inequality

• **Theorem:** Let $X_1, X_2, \dots, X_n \in \mathcal{X}$ be independent random variables and $f: \mathcal{X}^n \rightarrow \mathbb{R}$ be a function that is **stable** to input change:

$$\sup_{x_1, x_2, \dots, x_n, x'_i} |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i$$

for all $i \in [1, n]$. Then, for all $\epsilon > 0$:

finite-sample bound!

$$P(|f(X_1, \dots, X_n) - \mathbb{E}f(X_1, \dots, X_n)| \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}\right)$$

e) Generalization Bound: Finite Hypothesis Space

• **Theorem:** Let \mathcal{H} be a **finite** hypothesis space, $|\mathcal{H}| < \infty$, then for any $\delta > 0$, with probability at least $1 - \delta$,

$$\forall h \in \mathcal{H}, \quad \mathcal{E}(h) \leq \hat{\mathcal{E}}_{\mathcal{D}_n}(h) + \sqrt{\frac{\log|\mathcal{H}| + \log \frac{2}{\delta}}{2n}}$$

f) Rademacher Complexity

- Let \mathcal{G} be a family of general functions mapping from \mathcal{Z} to $[0,1]$.
- Let σ_i be **i.i.d.** uniform random variables (i.e. **Rademacher variable**):

$$P(\sigma_i = 1) = 1/2, P(\sigma_i = -1) = 1/2$$
- **Empirical Rademacher Complexity** of \mathcal{G} on a size- n sample set $\mathcal{S}_n = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\}$:

Average loss on all dichotomies

$$\hat{\mathcal{R}}_{\mathcal{S}_n}(\mathcal{G}) = \mathbb{E}_{\mathbf{g} \in \mathcal{G}} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{z}_i) \right]$$

Peter Bartl

$$\mathcal{R}_n(\mathcal{G}) = \mathbb{E}_{\mathcal{S}_n \sim D^n} \hat{\mathcal{R}}_{\mathcal{S}_n}(\mathcal{G}) = \mathbb{E}_{\mathcal{S}_n \sim D^n} \mathbb{E}_{\mathbf{g} \in \mathcal{G}} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(\mathbf{z}_i) \right]$$

- **Theorem:** Rademacher complexity will decrease when **n increases**:

$$\mathcal{R}_{n+1}(\mathcal{G}) \leq \mathcal{R}_n(\mathcal{G})$$
- **Theorem:** Let \mathcal{H} be a family of binary classifiers taking values in $\{-1, +1\}$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, the generalization bound holds for all $h \in \mathcal{H}$:

$$\mathcal{E}_{\mathcal{D}}(h) \leq \hat{\mathcal{E}}_{\mathcal{D}_n}(h) + \mathcal{R}_n(\mathcal{H}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

$$\mathcal{E}_{\mathcal{D}}(h) \leq \hat{\mathcal{E}}_{\mathcal{D}_n}(h) + \hat{\mathcal{R}}_n(\mathcal{H}) + 3 \sqrt{\frac{\log(2/\delta)}{2n}}$$

g) VC dim

- **VC-dimension** of a hypothesis space \mathcal{H} is defined by $\text{VCdim}(\mathcal{H}) = \max\{n: \Pi_{\mathcal{H}}(n) = 2^n\}$
 - VC-dimension is essentially the size of the **largest set** that can be fully shattered by \mathcal{H} .
 - No need to shatter every sample set of size n .
- **Lower bound:** Build a sample of size n that can be shattered.
- **Upper bound:** Prove all samples of size $n + 1$ **cannot be shattered**.
- Let \mathcal{H} be a hypothesis set with $\text{VCdim}(\mathcal{H}) = d$
 - If $d = \infty$, $\Pi_{\mathcal{H}}(n) \leq 2^n$.

$d < \infty$ is PAC learnable!
 - If $d < \infty$, $\Pi_{\mathcal{H}}(n) \leq \left(\frac{en}{d}\right)^d = O(n^d)$.
- **Theorem:** Let \mathcal{H} be a family of functions taking values in $\{-1, +1\}$ with VC-dimension d , then for any $\delta > 0$, with probability at least $1 - \delta$, for all $h \in \mathcal{H}$:

$$\mathcal{E}_{\mathcal{D}}(h) \leq \hat{\mathcal{E}}_{\mathcal{D}_n}(h) + \sqrt{\frac{2d \log \frac{en}{d}}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$
- The general form: $\mathcal{E}_{\mathcal{D}}(h) \leq \hat{\mathcal{E}}_{\mathcal{D}_n}(h) + O\left(\sqrt{\frac{\log(n/d)}{(n/d)}}.$

7. Decision Tree

- Misclassification error:

$$\text{Err}(\mathcal{D}) = 1 - \max_{1 \leq k \leq K} \left(\frac{|C_k|}{|\mathcal{D}|} \right)$$

- Entropy (used in ID3 and C4.5):

$$H(\mathcal{D}) = - \sum_{k=1}^K \frac{|C_k|}{|\mathcal{D}|} \log \frac{|C_k|}{|\mathcal{D}|}$$

- Gini index (used in CART):

$$\text{Gini}(\mathcal{D}) = 1 - \sum_{k=1}^K \left(\frac{|C_k|}{|\mathcal{D}|} \right)^2$$

a) ID3

Class label

ID3(Examples, Target.attribute, Attributes)

- create a **Root** node for the tree; assign all **Examples** to **Root**;

Stop Criteria
- if all **Examples** are positive, return the single-node tree **Root**, with label=+;
- if all **Examples** are negative, return the single-node tree **Root**, with label=-;
- if **Attributes** is empty, return the single-node tree **Root**, with label = the most common value of **Target.attribute** in **Examples**;
- otherwise // **Main loop**:

$A \leftarrow$ the attribute from **Attributes** that best* classifies **Examples**;
 the decision attribute for **Root** $\leftarrow A$;
 for each possible value v_i of A

$O(dn)$
in each layer

add a new tree branch below **Root**, corresponding to the test $A = v_i$;
 let Examples_{v_i} be the subset of **Examples** that have the value v_i for A ;
 if Examples_{v_i} is empty
 below this new branch add a leaf node with label = the most common value of **Target.attribute** in **Examples**;
 else
 below this new branch add the subtree
 ID3(Examples_{v_i} , **Target.attribute**, **Attributes** \ { A });

depth
 $\min(d, \log n)$
- return **Root**;

The best attribute is the one with the highest information gain.

b) C4.5

IG Rate

• By penalizing **multivalued rate**, IG is improved to **Gain Ratio** (GR):

$$\text{GR} = \frac{\text{Information Gain}}{\text{multivalued rate}}$$

• C4.5 Algorithm measures multivalued rate by **Intrinsic Value** (IV):

$$\text{IV}(F) = - \sum_{i=1}^{|V|} \frac{|\mathcal{D}_i|}{|\mathcal{D}|} \log \frac{|\mathcal{D}_i|}{|\mathcal{D}|}$$

- When \mathcal{D} is split by the feature f to \mathcal{D}_i if feature value is i .
- Which is the **entropy of the value probability** of the feature f .

Attribute with costs

• **C4.5** introduces a new criterion that could penalize the cost:

$$\frac{(\text{Gain Ratio})^2}{\text{Cost}}$$

Missing Value

- **Step 1:** When computing IG, for a feature with missing values:
 - Compute the **ratio of samples** with missing value: $|\bar{\mathcal{D}}|/|\mathcal{D}| = \rho$.
 - Compute the IG on $\mathcal{D} \setminus \bar{\mathcal{D}}$ (all samples with values):
 - The IG for **this feature** is computed by: $(1 - \rho) \cdot \text{IG}$.
 - Then compute **Gain Ratio** based on this new Information Gain.

Continuous Variables

- If x_{1j}, \dots, x_{nj} are the **sorted values** of the j th feature for n instances.
 - We only need to check split points **between adjacent values**
 - Traditionally take split points **halfway between adjacent values**:

$$s_j \in \left\{ \frac{1}{2} (x_{rj} + x_{(r+1)j}) \mid r = 1, \dots, n-1 \right\}$$
 - Check performance of $n - 1$ splits and find the one with highest IG.
- Pruning
- Pre-Pruning
 - If the accuracy grows no more than eta or even decreases: Stop the split
 - May cause underfitting
 - Post-Pruning
 - Avoid underfitting

• The **cost complexity criterion** with parameter α for tree T is

$$C_{\alpha}(T) = \hat{\mathcal{E}}(T) + \alpha |T| = \sum_{t=1}^{|T|} N_t H_t(T) + \alpha |T|$$

$$= - \sum_{t=1}^{|T|} \sum_{k=1}^K N_{tk} \log \frac{N_{tk}}{N_t} + \alpha |T|$$

- $|T|$: #leaf nodes, N_{tk} is the #examples of class k in leaf node t
- $\hat{\mathcal{E}}(T)$ is the empirical error of the tree on training data.
- **Cost Complexity Pruning:** (Given tree T_0 and parameter α)
 - Compute the empirical entropy $H_t(T)$ of each node
 - Recursively shrink from leaf nodes to internal nodes
 - If $C_{\alpha}(T_A) \leq C_{\alpha}(T_A)$: prune leaf A and use parent B as new leaf
- c) CART

Classification Tree

- If \mathcal{D} is split into \mathcal{D}_1 and \mathcal{D}_2 by whether attribute **$A(x) = a$** :

$$\mathcal{D}_1 = \{(x, y) \in \mathcal{D}: A(x) = a\}, \mathcal{D}_2 = \mathcal{D} - \mathcal{D}_1$$

• Then under the condition of attribute A , the Gini index of \mathcal{D} is

$$\text{Gini}(\mathcal{D}, A) = \frac{|\mathcal{D}_1|}{|\mathcal{D}|} \text{Gini}(\mathcal{D}_1) + \frac{|\mathcal{D}_2|}{|\mathcal{D}|} \text{Gini}(\mathcal{D}_2)$$

• Gini index **Gini(D, A)** is used to evaluate node split in **CART**, which is a **binary** tree where each node is split into $A(x) = a$ and $A(x) \neq a$

Regression Tree

- For each **splitting variable** j and **splitting point** s
 - Assume that the region R is split by j and s to R_1 and R_2 .

$$\hat{\nu}_1(j, s) = \text{avg}(y_i | x_i \in R_1(j, s))$$

$$\hat{\nu}_2(j, s) = \text{avg}(y_i | x_i \in R_2(j, s))$$

• Find j, s that minimize the loss

$$\ell(j, s) = \underbrace{\sum_{(x_i \in R_1(j, s))} (y_i - \hat{\nu}_1(j, s))^2}_{\text{L2 loss for samples in } R_1(j, s)} + \underbrace{\sum_{(x_i \in R_2(j, s))} (y_i - \hat{\nu}_2(j, s))^2}_{\text{L2 loss for samples in } R_2(j, s)}$$

• As we did in classification, **search on all (j, s) and find a best one.**

8. Random Forest

a) Bagging

• **Bagging** is a general technique to **reduce the variance** of an estimator:

- Draw B **bootstrap samples** $\mathcal{D}_1^b, \dots, \mathcal{D}_B^b$ from original data \mathcal{D}_n .
- Let $h_{\mathcal{D}_1^b}, \dots, h_{\mathcal{D}_B^b}$ be the prediction functions for each sample.
- The **bagged prediction function** is a **combination** of these functions:

$$h_{\text{bag}}(x) = \text{Combine}(h_{\mathcal{D}_1^b}(x), \dots, h_{\mathcal{D}_B^b}(x))$$

- How can we combine:

- Binary/multiclass predictions } **Vote**
 - Class probability predictions }
 - Prediction functions for regression } **Average**

$\text{Var}(h_{\text{bag}}) \leq \text{Var}(h_{\mathcal{D}_n})$
- For each training point x_i , let

$$S_i = \{b | \mathcal{D}_n^b \text{ does not contain } x_i\}$$
- The **OOB prediction** on x_i is

$$h_{\text{OOB}}(x_i) = \frac{1}{|S_i|} \sum_{b \in S_i} h_{\mathcal{D}_n^b}(x_i)$$
- b) Breiman Algorithm

Input: Data set $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$;
Feature subset size K .

Process:

1. $N \leftarrow$ create a tree node based on D ;
2. **If all instances in the same class then return N**
3. $\mathcal{F} \leftarrow$ the set of features that can be split further;
4. **If \mathcal{F} is empty then return N**
5. $\tilde{\mathcal{F}} \leftarrow$ select K features from \mathcal{F} randomly;
6. $N.f \leftarrow$ the feature which has the best split point in $\tilde{\mathcal{F}}$;
7. $N.p \leftarrow$ the best split point on $N.f$;
8. $D_l \leftarrow$ subset of D with values on $N.f$ smaller than $N.p$;
9. $D_r \leftarrow$ subset of D with values on $N.f$ no smaller than $N.p$;
10. $N_l \leftarrow$ call the process with parameters (D_l, K) ;
11. $N_r \leftarrow$ call the process with parameters (D_r, K) ; $K \approx \sqrt{d}$
12. **return N**

Output: A random decision tree

Randomized features

each tree

