

1.	Model Selection
a)	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 $\hat{\epsilon}_{\mathcal{D}}(h^{(i)})$
2.	KNN
a)	误差界估计
b)	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:
c)	Distance
•	Cosine Distance: $\cos(\mathbf{x}_i, \mathbf{x}_j) = \frac{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\ \mathbf{x}_i\ \ \mathbf{x}_j\ }$
•	Minkowski Distance: $D_p(\mathbf{x}_i, \mathbf{x}_j) = \sqrt[p]{\sum_{k=1}^d x_{ik} - x_{jk} ^p}$
•	Mahalanobis Distance: $D_M(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{M} (\mathbf{x}_i - \mathbf{x}_j)}$
d)	Weighted KNN
•	Distance-weighted Classification:
	$\hat{y} = \operatorname{argmax}_{c \in \mathcal{Y}} \sum_{x' \in \text{KNN}_c(x)} \frac{1}{D(x, x')^2}$
•	Distance-weighted Regression:
	$\hat{y} = \sum_{x' \in \text{KNN}(x)} \frac{1}{D(x, x')^2} y'$
e)	Nearest Centroid Classifier
•	Given dataset $\{\mathbf{x}_i, y_i\}_{i=1}^N$ with class label $y_i \in \mathcal{Y}$:
-	Compute per-class centroids $\mu_c = \frac{1}{ N_c } \sum_{i: y_i=c} \mathbf{x}_i$
-	Prediction function: $\hat{y} = \operatorname{argmin}_{c \in \mathcal{Y}} \ \mathbf{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).
f)	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 $n \gg 2^d$.
g)	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)

3.	Linear Regression
a)	Analytic Sol
•	Computing the gradient of $\hat{\epsilon}(\mathbf{w})$ w.r.t. \mathbf{w} and setting it to zero yields the optimal parameter \mathbf{w}^* :
	$\hat{\epsilon}(\mathbf{w}) = \ \mathbf{X}\mathbf{w} - \mathbf{y}\ ^2$ $\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$ $\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$ $\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ Normal Equation
b)	Optimization: SGD
•	In each iteration $t \leq T$: - Randomly sample a minibatch of $m \ll n$ points $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$
-	Set $J^t(\mathbf{w}^t) = \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}^t; \mathbf{x}_i, y_i)$
-	Compute gradient on minibatch: $\Delta^t = \nabla_{\mathbf{w}} J^t(\mathbf{w}^t)$
-	Update parameters with learning rate η : $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \Delta^t$
•	For linear regression: $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - 2\eta \mathbf{X}_m^T(\mathbf{X}_m \mathbf{w}^t - \mathbf{y})$ - Computational complexity $O(dmT)$. What if large dimension d ?
•	Theoretically, convergence rate is $O(1/\sqrt{T})$ under convex condition. c) Regularization
•	Norm-regularization will make the hypothesis smooth at any point.
4.	Linear Classification
a)	Logistic Regression
•	To map the output of $h(\mathbf{x})$ 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 $h(\mathbf{x}) = p(y=1 \mathbf{x})$ as the probability to label \mathbf{x} as $y=1$.
	$\hat{\epsilon}(\mathbf{w}) = -\sum_{i=1}^n \{y_i \log \sigma(h_{\mathbf{w}}(\mathbf{x}_i)) + (1 - y_i) \log[1 - \sigma(h_{\mathbf{w}}(\mathbf{x}_i))]\} + \lambda \sum_{j=1}^d w_j^2$
b)	SoftMax Regression
•	Softmax function normalizes multiple outputs in a probability vector:
	$p(y=c \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x})}$ Sum over all classes
	Loss for each data point (\mathbf{x}_i, y_i) :
	$\ell(h(\mathbf{x}_i), y_i) = \begin{cases} -\log \frac{\exp(\mathbf{w}_i^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = 1 \\ -\log \frac{\exp(\mathbf{w}_2^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = 2 \\ \vdots \\ -\log \frac{\exp(\mathbf{w}_C^T \mathbf{x}_i)}{\sum_{r=1}^C \exp(\mathbf{w}_r^T \mathbf{x}_i)}, & y_i = C \end{cases}$ Equivalent to maximum log-likelihood estimation
5.	SVM
a)	SM SVM
Original form	$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \ \mathbf{w}\ _2^2$ $\text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$ $\xi_i \geq 0, \sum_{i=1}^n \xi_i \leq n', 1 \leq i \leq n$
Lagrangian equivalent form	$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \ \mathbf{w}\ _2^2 + C \sum_{i=1}^n \xi_i$ $\text{s.t. } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$ $\xi_i \geq 0, 1 \leq i \leq n$
b)	Dual Problem

Primal Problem	$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ s.t. $g_j(\mathbf{x}) \leq 0 \text{ for } j = 1, \dots, J$ $h_k(\mathbf{x}) = 0 \text{ for } k = 1, \dots, K$										
•	With Lagrange multipliers λ, μ , the Lagrangian function is defined as										
	$L(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \sum_{j=1}^J \lambda_j g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k h_k(\mathbf{x})$										
KKT conditions	KKT conditions: for $1 \leq j \leq J$ - Primal feasibility: $g_j(\mathbf{x}) \leq 0, h_k(\mathbf{x}) = 0$ - Dual feasibility: $\lambda_j \geq 0$ - Complementary slackness: $\lambda_j g_j(\mathbf{x}) = 0$ - The gradient of the Lagrangian function w.r.t. \mathbf{x} vanishes to 0: $\nabla_{\mathbf{x}} f(\mathbf{x}) + \sum_{j=1}^J \lambda_j \nabla_{\mathbf{x}} g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k \nabla_{\mathbf{x}} h_k(\mathbf{x}) = 0$										
•	If $\mathbf{x}^*, \lambda^*, \mu^*$ satisfy KKT for a convex problem, then \mathbf{x}^* is an optimum: $f(\mathbf{x}^*) = L(\mathbf{x}^*, \lambda^*, \mu^*)$ - \mathbf{x}^* is optimal if and only if there exist λ, μ satisfying KKT conditions.										
Dual Problem	Primal objective as $\Pi(\mathbf{x}) = \max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} L(\mathbf{x}, \lambda, \mu)$										
•	Define the dual objective as $\Gamma(\lambda, \mu) \triangleq \min_{\mathbf{x} \in \mathbb{R}^d} L(\mathbf{x}, \lambda, \mu)$										
•	Weak duality: $\max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} \Gamma(\lambda, \mu) \leq \min_{\mathbf{x} \in \mathbb{R}^d} \Pi(\mathbf{x})$ Minmax theorem: $\max_x \min_y f(x, y) \leq \min_x \max_y f(x, y)$										
•	Dual Problem: $\max_{\lambda \in \mathbb{R}^J, \mu \in \mathbb{R}^K} \Gamma(\lambda, \mu)$ s.t. $\lambda_j \geq 0 \text{ for } j = 1, \dots, J$										
Soft-Margin SVM Dual Problem	Lagrangian function (with $2n$ inequality constraints): $L(\mathbf{w}, b, \alpha, \xi, \mu)$ $= \frac{1}{2} \ \mathbf{w}\ _2^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i(\mathbf{w} \cdot \mathbf{x}_i + b)) - \sum_{i=1}^n \mu_i \xi_i$ $\alpha_i \geq 0, \mu_i \geq 0, i = 1, \dots, n$										
•	After solving for α , we can solve for $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$. Dual Problem of Soft-SVM: $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$ s.t. $\sum_{i=1}^n \alpha_i y_i = 0$ $0 \leq \alpha_i \leq C, 1 \leq i \leq n$										
•	- After solving for α , we can solve for $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$. - Testing: $f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x}) + b = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$. - During training and testing, only need to compute $k(\mathbf{x}_i, \mathbf{x}_j)$.										
6.	Learning Theory										
a)	Bias-Variance Decomposition										
•	Theorem: $E_{\mathcal{D}}(\mathbf{x}) = \text{Var}(y \mathbf{x}) + \text{Bias}(h_{\mathcal{D}}(\mathbf{x}) \mathbf{x})^2 + \text{Var}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x}) \mathbf{x}]$										
SS											
b)	Approximation Error and Estimation Error										
•	General equality: given a target function f , for any $h \in \mathcal{H}$, $E(h) - E^*(f) = [\underbrace{E(h) - E(h^*)}_{\text{estimation}}] + [\underbrace{E(h^*) - E^*(f)}_{\text{approximation}}]$										
SS											
c)	PAC Learning										
•	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 $\mathbf{g}_i^t = \begin{cases} \mathbf{w}^t & \text{if } y_i(\mathbf{w}^t \cdot \mathbf{x}_i + b) \geq 1 \\ \mathbf{w}^t - C y_i \mathbf{x}_i & \text{otherwise} \end{cases}$ - Update $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta \mathbf{g}_i^t$ Output $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}^t$ Convergence rate: $O(\frac{1}{\sqrt{T}})$ (Eq. 14.3)										
•	Commonly used kernel functions for vector data:										
	<table border="1"><tr><td>Linear</td><td>$k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2$</td></tr><tr><td>Polynomial</td><td>$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^d$</td></tr><tr><td>Polynomial (with constant terms)</td><td>$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + 1)^d$</td></tr><tr><td>RBF (Gaussian kernel)</td><td>$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\ \mathbf{x}_1 - \mathbf{x}_2\ ^2}{2\sigma^2}\right)$</td></tr><tr><td>Hyperbolic Tangent</td><td>$k(\mathbf{x}_1, \mathbf{x}_2) = \tanh(\beta \mathbf{x}_1 \cdot \mathbf{x}_2 + c)$</td></tr></table>	Linear	$k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2$	Polynomial	$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^d$	Polynomial (with constant terms)	$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + 1)^d$	RBF (Gaussian kernel)	$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\ \mathbf{x}_1 - \mathbf{x}_2\ ^2}{2\sigma^2}\right)$	Hyperbolic Tangent	$k(\mathbf{x}_1, \mathbf{x}_2) = \tanh(\beta \mathbf{x}_1 \cdot \mathbf{x}_2 + c)$
Linear	$k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \cdot \mathbf{x}_2$										
Polynomial	$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2)^d$										
Polynomial (with constant terms)	$k(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1 \cdot \mathbf{x}_2 + 1)^d$										
RBF (Gaussian kernel)	$k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\ \mathbf{x}_1 - \mathbf{x}_2\ ^2}{2\sigma^2}\right)$										
Hyperbolic Tangent	$k(\mathbf{x}_1, \mathbf{x}_2) = \tanh(\beta \mathbf{x}_1 \cdot \mathbf{x}_2 + c)$										
•	From kernel functions k_1, k_2 , we can construct new kernel functions: - $k'(\mathbf{x}_1, \mathbf{x}_2) = k_1 \otimes k_2(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) k_2(\mathbf{x}_1, \mathbf{x}_2)$ - 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 $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, the kernel matrix is semi-definite ($\mathbf{K} \geq 0$, i.e. $\mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$).										
•	For kernel functions k_1, k_2, \dots, k_s , and $\gamma_1, \gamma_2, \dots, \gamma_s$, - $\sum_{i=1}^s \gamma_i k_s$ is also (multi-)kernel function, because $\sum_{i=1}^s \gamma_i k_s \geq 0$.										
•	Dual Problem of Soft-SVM: $\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$ s.t. $\sum_{i=1}^n \alpha_i y_i = 0$ $0 \leq \alpha_i \leq C, 1 \leq i \leq n$										
SS											
•	- After solving for α , we can solve for $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$. - Testing: $f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x}) + b = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$.										
b)	Approximation Error and Estimation Error										
•	General equality: given a target function f , for any $h \in \mathcal{H}$, $E(h) - E^*(f) = [\underbrace{E(h) - E(h^*)}_{\text{estimation}}] + [\underbrace{E(h^*) - E^*(f)}_{\text{approximation}}]$										
c)	PAC Learning										
•	$\min_{\mathbf{w}, b} \frac{1}{2} \ \mathbf{w}\ _2^2 + \frac{C}{n} \sum_{i=1}^n \max[0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b)]$										
•	$\begin{aligned} & \bullet \quad \ell_{0/1} = 1[t \leq 0] \\ & \bullet \quad r_{\text{hinge}} = \max[1 - t, 0] \\ & \bullet \quad r_{\text{exp}} = \exp(-t) \\ & \bullet \quad \text{logistic} = \log(1 + e^{-t}) \\ & \bullet \quad \text{where } t = \mathbf{y}(\mathbf{w} \cdot \mathbf{x}) + b \end{aligned}$										

A hypothesis space \mathcal{H} is **PAC-learnable** if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}()$, 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** $n \geq \text{poly}\left(\frac{1}{\epsilon}, \frac{1}{\delta}, |\mathcal{H}|\right)$:

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



Basic

• Union bound: $P(A \vee B) \leq P(A) + P(B)$.

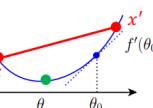
• Inversion: if $P(X \geq \epsilon) \leq f(\epsilon)$, then for any $\delta > 0$, with probability at least $1 - \delta$, $X \leq f^{-1}(\delta)$.

- Useful to interchange ϵ and δ : $P(|\hat{\mathcal{E}}_{D_n}(h) - \mathcal{E}(h)| \geq \epsilon) \leq \delta$.

• Jensen's inequality: if f is convex, then $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$.

- Convex function: f is convex if $\forall \lambda \in [0, 1], x, x' \in \text{dom}(f)$:

$$\begin{aligned} f(\lambda x + (1 - \lambda)x') \\ \leq \lambda f(x) + (1 - \lambda)f(x') \\ - f(x') \geq f(x) + g^T(x' - x) \\ \bullet g: \text{Subgradient.} \end{aligned}$$



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:

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

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

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

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 $1 - \delta$,

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

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 $S_n = \{z_1, z_2, \dots, z_n\}$: Average loss on all dichotomies

$$\hat{\mathcal{R}}_{S_n}(\mathcal{G}) = \mathbb{E}_{\sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right]$$



• (Expected) Rademacher Complexity:

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

• **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 $h \in \mathcal{H}$:

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

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

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 $d = \infty, \Pi_{\mathcal{H}}(n) \leq 2^n$.



- 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 $h \in \mathcal{H}$:

$$\mathcal{E}_D(h) \leq \hat{\mathcal{E}}_{D_n}(h) + \sqrt{\frac{2d \log \frac{en}{d}}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

• The general form: $\mathcal{E}_D(h) \leq \hat{\mathcal{E}}_{D_n}(h) + O\left(\sqrt{\frac{\log(n/d)}{n/d}}\right)$.

Decision Tree

- Misclassification error:

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

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

$$H(\mathcal{D}) = - \sum_{k=1}^K \frac{|\mathcal{C}_k|}{|\mathcal{D}|} \log \frac{|\mathcal{C}_k|}{|\mathcal{D}|} \quad \boxed{E}$$

- Gini index (used in CART):

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

a) ID3

Class label

ID3(Examples, Target_attribute, Attributes)

Stop Criteria

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

 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});
 * 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 I: When computing IG, for a feature with missing values:

- Compute the ratio of samples with missing value: $|\mathcal{D}_m|/|\mathcal{D}| = \rho$.

- Compute the IG on $\mathcal{D} \setminus \mathcal{D}_m$ (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_1, \dots, x_n 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

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

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

- $|T|$: #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 $\mathcal{C}_\alpha(T_B) \leq \mathcal{C}_\alpha(T_A)$: 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

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

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

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

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

• Find j, s that minimize the loss

$$\ell(j, s) = \sum_{x_i \in R_1(j, s)} (y_i - \hat{v}_1(j, s))^2 + \sum_{x_i \in R_2(j, s)} (y_i - \hat{v}_2(j, s))^2$$

L2 loss for samples in $R_1(j, s)$ L2 loss for samples in $R_2(j, s)$

• As we did in classification, search on all (j, s) '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 $\mathcal{D}_1^B, \dots, \mathcal{D}_n^B$ from original data \mathcal{D}_n .

• Let $h_{D_1}^B, \dots, h_{D_n}^B$ be the prediction functions for each sample.

• The bagged prediction function is a combination of these functions:

$$h_{\text{bag}}(x) = \text{Combine}(h_{D_1}^B(x), \dots, h_{D_n}^B(x))$$

• How can we combine:

{ Binary/multiclass predictions } Vote $\text{Var}(h_{\text{bag}})$ $\text{Var}(h_{D_i})$

{ Class probability predictions }

{ Prediction functions for regression } Average

• For each training point x_i , let

$$S_i = \{b \mid \mathcal{D}_i^B \text{ does not contain } x_i\}$$

• The OOB prediction on x_i is

$$h_{\text{OOB}}(x_i) = \frac{1}{|S_i|} \sum_{b \in S_i} h_{D_b}^B(x_i)$$

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. $F \leftarrow$ the set of features that can be split further;

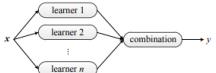
4. **if F is empty then return N**

5. $\tilde{F} \leftarrow$ select K features from F randomly;

6. $N_f \leftarrow$ the feature which has the best split point in \tilde{F} ;

9. Boosting

a) Ensemble



• Parallel ensemble

- Each model is **built independently**

- Bagging and Random Forest

- Idea: Combine different strong models to **avoid overfitting**

• Sequential ensemble

- Models are **built sequentially**

- Boosting and XGBoost

- Idea: Combine many weak models to **avoid underfitting**.

- Try to add new models that do well where previous models lack

b) Ada Boosting



$H \subseteq \{-1, +1\}^X$.

ADABoost($S = ((x_1, y_1), \dots, (x_m, y_m))$, $m = n$)

```

1 for i ← 1 to m do
2    $D_i(i) \leftarrow \frac{1}{m}$ 
3 for t ← 1 to T do
4    $h_t \leftarrow$  base classifier in  $H$  with small error  $\epsilon_t = \Pr_{i \sim D_t}[h_t(x_i) \neq y_i]$ 
5    $\alpha_t \leftarrow \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$ 
6    $Z_t \leftarrow 2[\epsilon_t(1 - \epsilon_t)]^{\frac{1}{2}}$  ▷ normalization factor
7   for i ← 1 to m do
8      $D_{t+1}(i) \leftarrow \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$  (Weight of instance i at iteration t + 1)
9    $f_t \leftarrow \sum_{s=1}^t \alpha_s h_s$ 
10 return  $h = \text{sgn}(f_T)$ 
    
```

Why?



Boosting > Bagging > Tree

• Proof: We minimize $F(\bar{a})$ by Coordinate Descent (CD).

• Note that $F(\bar{a}_{t-1} + \eta e_k) = \frac{1}{n} \sum_{i=1}^n e^{-y_i \sum_{j=1}^t \bar{a}_{t-1,j} h_j(x_i)} - \eta y_i h_k(x_i)$

• Direction: for each unit vector e_k with best directional derivative:

$$F'(\bar{a}_{t-1}, e_k) = \lim_{\eta \rightarrow 0} \frac{F(\bar{a}_{t-1} + \eta e_k) - F(\bar{a}_{t-1})}{\eta}$$

$$= -\frac{1}{n} \sum_{i=1}^n y_i h_k(x_i) e^{-y_i \sum_{j=1}^t \bar{a}_{t-1,j} h_j(x_i)}$$

$$= -\frac{1}{n} \sum_{i=1}^n y_i h_k(x_i) \bar{D}_t(i) \bar{Z}_t$$

$$= -\left[\sum_{i=1}^n \bar{D}_t(i) \mathbf{1}_{[y_i h_k(x_i) = +1]} - \sum_{i=1}^n \bar{D}_t(i) \mathbf{1}_{[y_i h_k(x_i) = -1]} \right] \bar{Z}_t$$

$$= -[(1 - \bar{\epsilon}_{t,k}) - \bar{\epsilon}_{t,k}] \bar{Z}_t = [2\bar{\epsilon}_{t,k} - 1] \bar{Z}_t$$

Maximum descent direction k is the one minimizing $\bar{\epsilon}_{t,k}$ over distribution \bar{D}_t

• Step size: η is chosen to minimize $F(\bar{a}_{t-1} + \eta e_k)$:

$$\frac{\partial F(\bar{a}_{t-1} + \eta e_k)}{\partial \eta} = 0 \Rightarrow -\frac{1}{n} \sum_{i=1}^n y_i h_k(x_i) e^{-y_i \sum_{j=1}^t \bar{a}_{t-1,j} h_j(x_i)} - \eta y_i h_k(x_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^n y_i h_k(x_i) \bar{D}_t(i) \bar{Z}_t e^{-\eta y_i h_k(x_i)} = 0$$

$$\Leftrightarrow \sum_{i=1}^n \bar{D}_t(i) y_i h_k(x_i) e^{-\eta y_i h_k(x_i)} = 0$$

$$\Leftrightarrow \sum_{i=1}^n \bar{D}_t(i) \mathbf{1}_{[y_i h_k(x_i) = +1]} e^{-\eta} - \sum_{i=1}^n \bar{D}_t(i) \mathbf{1}_{[y_i h_k(x_i) = -1]} e^{\eta} = 0$$

$$\Leftrightarrow [(1 - \bar{\epsilon}_{t,k}) - \bar{\epsilon}_{t,k} e^{\eta}] = 0$$

$$\Leftrightarrow \eta = \frac{1}{2} \log \frac{1 - \bar{\epsilon}_{t,k}}{\bar{\epsilon}_{t,k}}$$

Step size η of CD is weight α_t to classifier h_t in the t th iteration

• Base learners:

- Decision trees, quite often just decision stumps.

- Stumps are simple, **weak learners**:

• Think the XOR example.

• Boosting stumps:

- Data in \mathbb{R}^d .

- Associate a stump to each feature.

- Pre-sort each feature: $O(dn \log n)$

- At each round, find best feature and threshold.

- Total complexity with T stumps: $O(dn \log n + ndT)$

• Strengths:

- Simple: straightforward implementation.
- Efficient: complexity $O(ndT)$ for stumps.
- Theoretical guarantees: aim to maximize margin but not achieved.

• Outliers:

- AdaBoost assigns larger weights to **harder examples**.
- Detecting **mislabeled examples** and dealing with **noisy data**.

• Weaknesses:

- Determining T , the number of rounds of boosting is **nontrivial**.
- Need to determine base learners: risk of **overfitting** or **low margin**.

• Theorem: The empirical error of the classifier h output by AdaBoost verifies:

$$\hat{E}(h) \leq \exp \left[-2 \sum_{t=1}^T \left(\frac{1}{2} - \epsilon_t \right)^2 \right]$$

Furthermore, if for all $t \in [T]$, $\gamma \leq \left(\frac{1}{2} - \epsilon_t \right)$, then

$$\hat{E}(h) \leq \exp[-2\gamma^2 T]$$

• Additive model $f(x) = \sum_{t=1}^T \alpha_t h_t(x)$ and $h(x) = \text{sgn}(f(x))$.

- ϵ_t is the error of each base learner h_t .

• Note: γ does not need to be known in advance: Adaptive Boosting.

• Proof: As we defined and derived previously, $D_{T+1}(i) = \frac{1}{n} e^{-y_i f(x_i)}$.

$$\begin{aligned} \hat{E}(h) &\leq \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{[y_i f(x_i) \leq 0]} \leq \frac{1}{n} \sum_{i=1}^n e^{-y_i f(x_i)} \\ &\leq \frac{1}{n} \sum_{i=1}^n \left[n \prod_{t=1}^T Z_t \right] D_{T+1}(i) = \prod_{t=1}^T Z_t \end{aligned}$$

• Now recall that Z_t is a normalization factor:

$$\begin{aligned} Z_t &= \sum_{i=1}^n D_t(i) e^{-\alpha_t y_i h_t(x_i)} \\ &= \sum_{i: y_i h_t(x_i) \geq 0} D_t(i) e^{-\alpha_t} + \sum_{i: y_i h_t(x_i) < 0} D_t(i) e^{\alpha_t} \\ &= (1 - \bar{\epsilon}_t) e^{-\alpha_t} + \bar{\epsilon}_t e^{\alpha_t} \end{aligned}$$

• Proof: We have $\hat{E}(h) \leq \prod_{t=1}^T Z_t$ and $Z_t = (1 - \bar{\epsilon}_t) e^{-\alpha_t} + \bar{\epsilon}_t e^{\alpha_t}$.

• Derivative: $\frac{\partial \hat{E}}{\partial \alpha_t} = -(1 - \bar{\epsilon}_t) e^{-\alpha_t} + \bar{\epsilon}_t e^{\alpha_t} = 0 \Rightarrow \alpha_t = \frac{1}{2} \log \frac{1 - \bar{\epsilon}_t}{\bar{\epsilon}_t}$.

• Compute minimum: $Z_t = (1 - \bar{\epsilon}_t) e^{-\alpha_t} + \bar{\epsilon}_t e^{\alpha_t} = 2\sqrt{\bar{\epsilon}_t(1 - \bar{\epsilon}_t)}$.

$$\begin{aligned} \hat{E}(h) &\leq \prod_{t=1}^T Z_t \min_{t=1}^T 2\sqrt{\bar{\epsilon}_t(1 - \bar{\epsilon}_t)} \\ &= \prod_{t=1}^T \sqrt{1 - 4 \left(\frac{1}{2} - \bar{\epsilon}_t \right)^2} \quad [1 - x \leq e^{-x}] \\ &\leq \prod_{t=1}^T \exp \left[-2 \left(\frac{1}{2} - \bar{\epsilon}_t \right)^2 \right] = \exp \left[-2 \sum_{t=1}^T \left(\frac{1}{2} - \bar{\epsilon}_t \right)^2 \right] \end{aligned}$$

c) Gradient Boosting

• Gradient Boosting:

- For $u_{ti} = \alpha_t h_t(x_i)$, let the t th base learner $h_t(x_i)$ approximate

$$h_t(x_i) \approx -\nabla_{v_{ti}} \sum_{i=1}^n \ell(y_i, v_{ti}) = -\nabla_{v_{ti}} \ell(y_i, v_{ti})$$

- And α_t be the **learning rate** η .

• Gradient boosting applies **whenever**:

- Loss function is **sub-differentiable** w.r.t. training predictions $f(x_i)$.

- Can do **regression** on base hypotheses $h_t(x_i)$, e.g. regression tree.

1. $\forall f_0(x) = 0_s$

2. For $t=1$ to T :

- 计算在各个数据点上的梯度 $\mathbf{g}_t = \left(\frac{\partial}{\partial y_i} \ell(y_i, \hat{y}_i) \Big|_{\hat{y}_i=f_{t-1}(x_i)} \right)_{i=1}^n$
- 根据 $-\mathbf{g}_t$ 拟合一个回归模型, $h_t = \arg \min_{h \in \mathcal{F}} \sum_{i=1}^n (h(x_i) + g_{t,i})^2$
- 选择合适的步长 α_t , 最简单的选择是固定步长 $\eta \in (0, 1]$ 。
- 更新模型, $f_t(x) = f_{t-1}(x) + \alpha_t h_t(x)$

L2 Boosting

$$h_t = \arg \min_h \sum_{i=1}^n (h(x_i) - 2(y_i - f_{t-1}(x_i)))^2$$

Binomial Boosting

$$h_t = \arg \min_h \sum_{i=1}^n \left(h(x_i) - \frac{y_i}{1 + e^{y_i f_{t-1}(x_i)}} \right)^2$$

10. Reinforcement Learning: Model-Based

a) Mathematical Framework

Bellman Equation

• Similarly decomposed into **immediate reward** + **discounted value**.

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \\ q_\pi(s, a) &= \mathbb{E}_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a] \end{aligned}$$

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a) \quad q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$

Optimal

• Definition:

The **optimal state value function** v_* is the maximum state value function over all policies

$$v_*(s) = \max_\pi v_\pi(s)$$

The **optimal action value function** q_* is the maximum action value function over all policies

$$q_*(s, a) = \max_\pi q_\pi(s, a)$$

• Theorem: For any Markov Decision Process (MDP):

- There **exists** an optimal policy π_* that is better than or equal to all other policies:

$$\pi_* \geq \pi, \forall \pi$$

- All optimal policies achieve the optimal state value function:

$$v_{\pi_*}(s) = v_*(s)$$

- All optimal policies achieve the optimal action value function:

$$q_{\pi_*}(s, a) = q_*(s, a)$$

$$v_*(s) = \max_a q_*(s, a) \quad \text{Proof next slide}$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

$$v_*(s) = \max_a R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \max_{a'} q_*(s', a')$$

b) Policy Iteration

Evaluation • Solution: iterative application of Bellman expectation backup

- At each iteration $k+1$, for all states $s \in \mathcal{S}$:

- Update $v_{k+1}(s)$ from $v_k(s')$, where s' is a **successor** state of s

$$v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

Improvement

• Consider a **deterministic policy** $a = \pi(s)$.

• We can improve the policy by acting **greedily**

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} q_\pi(s, a)$$

c) Value Iteration

Converge to v_* by
Contraction Mapping

11. Reinforcement Learning: Model-Free
a) Monte Carlo Method

• First-visit MC prediction, for policy evaluation with $V \approx v_\pi$.

• The first time-step t that state s is visited in an episode:

- Increment counter $N(s) \leftarrow N(s) + 1$

- Increment total return $G(s) \leftarrow G(s) + R_t$

• Value is estimated by sample average return $V(s) = \frac{G(s)}{N(s)}$

• By law of large numbers: $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

• Every-visit MC update (online): $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$

$$\begin{aligned} \mu_k &= \frac{1}{k} \sum_{j=1}^k x_j = \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right) = \frac{1}{k} (x_k + (k-1)\mu_{k-1}) \\ &= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1}) = \mu_{k-1} + \alpha(x_k - \mu_{k-1}) \end{aligned}$$

• First-visit MC control (for ϵ -greedy policy), estimates $\pi \approx \pi_*$

• Initialize:

- $\pi \leftarrow$ an ϵ -soft policy, value $Q(s, a)$, empty list $\text{Returns}(s, a), \forall s, a$

• Repeat (for each episode):

- Sample an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

- $G \leftarrow 0$

• Repeat (for each step of episode, $t = T-1, T-2, \dots, 0$):

- $G \leftarrow \gamma G + R_{t+1}$

- Unless the pair S_t, A_t appears in $S_0, A_0, \dots, S_{t-1}, A_{t-1}$:

- Append G to $\text{Returns}(S_t, A_t)$

- $Q(S_t, A_t) \leftarrow \text{Avg}[\text{Returns}(S_t, A_t)]$

- $A^* \leftarrow \arg \max_a Q(S_t, A_t)$ Generalized policy iteration (GPI)

- For all $a \in \mathcal{A}$: $\pi(a|S_t) = \begin{cases} 1 - \epsilon + \epsilon/J, & \text{if } a = A^* \\ \epsilon/J, & \text{if } a \neq A^* \end{cases}$

• Theorem (ϵ -greedy policy improvement): For any ϵ -greedy policy π , the ϵ -greedy policy π' wrt q_π is an improvement, $v_{\pi'}(s) \geq v_\pi(s)$.

• Proof:

$$\begin{aligned} q_\pi(s, \pi'(s)) &= \sum_{a \in \mathcal{A}} \pi'(a|s) q_\pi(s, a) \\ &= \frac{\epsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \max_{a \in \mathcal{A}} q_\pi(s, a) \\ &\geq \frac{\epsilon}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} q_\pi(s, a) + (1 - \epsilon) \sum_{a \in \mathcal{A}} \frac{\pi(a|s) - \epsilon}{|\mathcal{A}|} q_\pi(s, a) \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a) = v_\pi(s) \end{aligned}$$

• By policy improvement theorem: $q_\pi(s, \pi'(s)) \geq v_\pi(s) \Rightarrow v_{\pi'}(s) \geq v_\pi(s)$.

b) TD Learning

- TD(0) prediction, for estimating $V \approx v_\pi$
- Input: the policy π to be evaluated, step size $\alpha \in [0,1]$
- Initialize:
 - Value $V(s), \forall s \in \mathcal{S}$, except that $V(\text{terminal}) = 0$
- Repeat (for each episode):
 - Initialize S
 - Repeat (for each step of episode):
 - $A \leftarrow$ action given by π for S
 - Take action A , observe R, S'
 - $V(S) \leftarrow V(S) + \alpha(R + \gamma V(S') - V(S))$
 - $S \leftarrow S'$
 - Until S is terminal
- N-step TD prediction, for estimateing $V \approx v_\pi$
- Initialize: Value $V(s), \forall s \in \mathcal{S}$
- Repeat (for each episode):
 - Initialize and store $S_0 \neq \text{terminal}, T \leftarrow \infty$ (big cost)
 - Repeat (for $t = 0, 1, 2, \dots$):
 - If $t < T$:
 - Take an action according to $\pi(\cdot | S_t)$
 - Observe and store next reward R_{t+1} & next state S_{t+1}
 - If S_{t+1} is terminal, then $T \leftarrow t + 1$
 - If $t \geq 0$:
 - $G \leftarrow \sum_{i=t+1}^{\min(\tau+n,T)} \gamma^{i-t-1} R_i$
 - If $\tau + n < T$, then: $G \leftarrow G + \gamma^n V(S_{\tau+n})$
 - $V(S_t) \leftarrow V(S_t) + \alpha(G - V(S_t))$
 - Until $\tau = T - 1$

- TD(1) lambda
- Input: the policy π , discount γ , decay rate λ .
 - Repeat (for each episode)
 - $E(s) \leftarrow 0$ for all $s \in \mathcal{S}$
 - Initialize S
 - Repeat (for each step of episode)
 - Choose $A \sim \pi(\cdot | S)$
 - Take action A , observe R, S'
 - TD error: $\delta = R + \gamma V(S') - V(S)$
 - Eligibility trace: $E(S) \leftarrow E(S) + 1$
 - For all $s \in \mathcal{S}$
 - $V(s) \leftarrow V(s) + \alpha \delta E(s)$
 - Eligibility trace: $E(s) = \gamma \lambda E(s)$
 - $S \leftarrow S'$
 - Until S is terminal

- Advantages of Temporal-Difference (TD) over Monte Carlo (MC):
 - Bootstrapping
 - Online, fully incremental updates, efficient
 - Applicable to continuing tasks without episodes

c) Sarsa

- Updating action-value function with Sarsa (S, A, R, S', A'):

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$$
- Initialize: $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ arbitrarily, $Q(\text{terminal}, \cdot) = 0$
- Repeat (for each episode):
 - Initialize S
 - Choose A from S using policy derived from Q (e.g. ϵ -greedy)
 - Repeat (for each step of episode):
 - Take action A , observe R, S' On-policy
 - Choose A' from S' using policy derived from Q (e.g. ϵ -greedy)
 - $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
 - $S \leftarrow S', A \leftarrow A'$
 - Until S is terminal

Sarsa (lambda)

- Initialize: $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
- Repeat (for each episode):
 - $E(s, a) \leftarrow 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
 - Initialize S, A
 - Repeat (for each step of episode):
 - Take action A , observe R, S'
 - Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)
 - TD error: $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$
 - Eligibility trace: $E(S, A) \leftarrow E(S, A) + 1$
 - For all $s \in \mathcal{S}, a \in \mathcal{A}(s)$:
 - $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$
 - Eligibility trace: $E(s, a) \leftarrow \gamma \lambda E(s, a)$
 - $S \leftarrow S', A \leftarrow A'$
 - Until S is terminal
 - d) Q-Learning

-
- Initialize: $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ arbitrarily, $Q(\text{terminal}, \cdot) = 0$
 - Repeat (for each episode):
 - Initialize S
 - Repeat (for each step of episode):
 - Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 - Take action A , observe R, S'
 - $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
 - $S \leftarrow S'$
 - Until S is terminal

Q-learning : try to evaluate value function for greedy policy while following ϵ -greedy policy, therefore off-policy

Sarsa : try to evaluate value function for ϵ -greedy policy while following ϵ -greedy policy, therefore on-policy

e) Policy Gradient Method

- Value-based reinforcement learning
 - Approximate value or action-value function by parameter \mathbf{w}
 - $\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$ or $\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$
 - Generate policy from value function by greedy or ϵ -greedy
- Policy-based reinforcement learning
 - Directly parametrize the policy by parameter θ

$$\pi(a|s, \theta) = \mathbb{P}[A_t = a | S_t = s]$$
Exploration!

- Advantages:
 - Better convergence, stochastic policies $\pi(a|s, \theta) \in (0, 1)$
 - Effective in high-dimensional or continuous action spaces
- Disadvantages: inefficient and of high variance (can be mitigated)

- Theorem (policy gradient): For any differentiable policy $\pi_\theta(a|s)$ and $J(\theta)$, policy gradient does not involve derivative of state distribution:

$$\nabla J(\theta) \propto \mathbb{E}_\pi[\nabla_\theta \log \pi_\theta(S_t, A_t) \nabla \log \pi(A_t | S_t, \theta)]$$

f) REINFORCE

• REINFORCE: Monte Carlo Policy Gradient

- Input:
 - A differentiable policy parameterization $\pi(a|s, \theta)$
 - Algorithm parameter: step size $\alpha > 0$
- Initialize: policy parameter $\theta \in \mathbb{R}^d$ (e.g. to $\mathbf{0}$)
- Repeat (for each episode)
 - Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T \sim \pi(\cdot | \cdot, \theta)$
 - Repeat (for each step of the episode $t = 0, 1, \dots, T - 1$)
 - $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$
 - $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \log \pi(A_t | S_t, \theta)$ Eligibility Vector

• REINFORCE with Baseline: Monte Carlo Policy Gradient

- Input: policy parameterization $\pi(a|s, \theta)$, state-value function $\hat{v}(s, \mathbf{w})$
 - Algorithm parameters: step sizes $\alpha^\theta > 0, \alpha^\mathbf{w} > 0$
- Initialize: policy parameter $\theta \in \mathbb{R}^d$, state-value weight $\mathbf{w} \in \mathbb{R}^d$
- Repeat (for each episode)
 - Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T \sim \pi(\cdot | \cdot, \theta)$
 - Repeat (for each step of the episode $t = 0, 1, \dots, T - 1$)
 - $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$
 - $\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$
 - $\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla \hat{v}(S_t, \mathbf{w})$
 - $\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla \log \pi(A_t | S_t, \theta)$

g) Actor-Critic

• QAC: Simple actor-critic algorithm based on action-value critic

- Input: policy parameterization $\pi(a|s, \theta)$, action-value critic $\hat{q}(s, a, \mathbf{w})$
 - Algorithm parameters: step sizes $\alpha^\theta > 0, \alpha^\mathbf{w} > 0$

- Initialize: policy parameter $\theta \in \mathbb{R}^d$, action-value weight $\mathbf{w} \in \mathbb{R}^d$

- Repeat (for each episode)
 - Initialize S
 - Repeat (for each step of the episode while S is not terminal)
 - $A \sim \pi(\cdot | S, \theta)$
 - Take action A , observe S', R , sample A'
 - $\delta \leftarrow R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$
 - $\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla \hat{q}(S, A, \mathbf{w})$
 - $\theta \leftarrow \theta + \alpha^\theta \hat{q}(S, A, \mathbf{w}) \nabla \log \pi(A | S, \theta)$
 - $S \leftarrow S', A \leftarrow A'$

- Input: policy parameterization $\pi(a|s, \theta)$, action-value critic $\hat{q}(s, a, \mathbf{w})$
 - Algorithm parameters: step sizes $\alpha^\theta > 0, \alpha^\mathbf{w} > 0$

- Initialize: policy parameter $\theta \in \mathbb{R}^d$, action-value weight $\mathbf{w} \in \mathbb{R}^d$

- Repeat (for each episode)
 - Initialize S
 - $e^\theta \leftarrow \mathbf{0}, e^\mathbf{w} \leftarrow \mathbf{0}$ (eligibility traces)
 - Repeat (for each step of the episode while S is not terminal)
 - $A \sim \pi(\cdot | S, \theta)$
 - Take action A , observe S', R , sample A'
 - $\delta \leftarrow R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$
 - $e^\mathbf{w} \leftarrow \gamma \lambda e^\mathbf{w} + \nabla \hat{q}(S, A, \mathbf{w})$
 - $e^\theta \leftarrow \gamma \lambda e^\theta + \nabla \log \pi(A | S, \theta)$
 - $\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta e^\mathbf{w}$
 - $\theta \leftarrow \theta + \alpha^\theta \hat{q}(S, A, \mathbf{w}) e^\theta$
 - $S \leftarrow S', A \leftarrow A'$

TD(λ)

Policy Gradient

- Policy gradient methods can follow any unbiased ascent direction

$$\text{Advantage function: } A_\pi(S_t, A_t) = q_\pi(S_t, A_t) - v_\pi(S_t)$$

- TD error: $\delta = R_{t+1} + \gamma v_\pi(S_{t+1}) - v_\pi(S_t)$

$$\nabla_\theta J(\theta) \propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(S_t, A_t) G_t]$$

REINFORCE

$$\propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(S_t, A_t) \hat{q}(S_t, A_t, \mathbf{w})]$$

Actor-Critic

$$\propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(S_t, A_t) \hat{A}(S_t, A_t, \mathbf{w})]$$

Advantage Actor-Critic

$$\propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(S_t, A_t) \delta]$$

TD Actor-Critic

$$\propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(S_t, A_t) \delta e]$$

TD(λ) Actor-Critic

- Use policy evaluation to update the critic (MC or TD prediction)