

Basis

Inner Product: $\langle \vec{x}, \vec{z} \rangle = \vec{x}^T \vec{z} = \sum_{i=1}^d x_i z_i$

Euclidean Norm: $\|\vec{x}\| = \sqrt{\langle \vec{x}, \vec{x} \rangle}$

Hadamard Product: $A \odot B = \begin{pmatrix} a_{11}b_{11} & \dots & a_{1n}b_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1}b_{m1} & \dots & a_{mn}b_{mn} \end{pmatrix}$

Kronecker Product: $A \otimes B = \begin{pmatrix} a_{11}B & \dots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \dots & a_{mn}B \end{pmatrix}$

Positive Semi-definite Matrix A

$\Leftrightarrow \forall \vec{v} \in \mathbb{R}^n: \vec{v}^T A \vec{v} \geq 0$

$\Leftrightarrow \forall \lambda_i: \lambda_i \geq 0$

Properties:

1. λ_i of $A \otimes B$ are the product of λ_i of A & B
2. $A \otimes B$ is a principal submatrix of $A \odot B$
3. $A \otimes B$ is positive semi-definite $\Leftrightarrow A \otimes B$ is as well

Moore-Penrose Pseudoinverse: $A^+ = (A^T A)^{-1} A^T$

Inner Product Space: $\langle \cdot, \cdot \rangle: \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$

1. $\langle x, y \rangle = \langle y, x \rangle$
2. $\langle x, x \rangle \geq 0$
3. $\langle \lambda_1 x + \lambda_2 y, z \rangle = \lambda_1 \langle x, z \rangle + \lambda_2 \langle y, z \rangle$

Computational Learning Theory

Risk: $R[f] = \int_{\mathcal{X} \times \mathcal{Y}} V(f(x), y) dP(x, y)$

Empirical Risk: $R_{\text{emp}}[f] = \frac{1}{n} \sum_{i=1}^n V(f(x_i), y_i)$

Regularization: Limit capacity

Lemma:

$$R[f] \leq R_{\text{emp}}[f] + \sqrt{\frac{VC_{\text{dim}}(\mathcal{H}) \left(\ln \frac{2m}{VC_{\text{dim}}(\mathcal{H})} + 1 \right) - \ln \frac{\delta}{2}}{m}}$$

Tikhonov Regularization:

$$\arg \min_{f \in \mathcal{H}} R_{\text{emp}}[f] + \lambda \Omega[f]$$

Version space: $VS_{\mathcal{H}, E} = \{h \in \mathcal{H}: h \text{ is consistent with } h\}$

General Boundary G : $g \in G \Leftrightarrow g$ is consistent $\wedge \exists g' \in \mathcal{H}: g' < g \wedge g'$ is consistent

Specific Boundary S : $s \in S \Leftrightarrow s$ is consistent $\wedge \exists s' \in \mathcal{H}: s' > s \wedge s'$ is consistent

Consistency:

1. h is complete: $E^+ \subseteq h$
2. h is correct: $E \cap h = \emptyset$

Generality: $h_1 \leq h_2: h_1$ is more general than h_2

PAC-learnable: $P(\text{error}(h) \leq \epsilon) \geq 1 - \delta$

Efficient PAC-learnable: Alg. runs in $\frac{1}{\epsilon}, \frac{1}{\delta}, n, \text{size}(C)$

Finite concept classes:

Theorem: C is PAC-learnable, if we draw m samples and find a consistent hypothesis with:

$$m \geq \frac{1}{\epsilon} (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

Corollary:

1. Theorem holds
2. m is polynomial in n
3. Alg. runs polynomial in m and n

$\Rightarrow C$ is efficiently PAC-learnable

Lemma: $VC_{\text{dim}}(C) \leq \log_2 |C|$

Results:

1. $k\text{-CNF}_n / k\text{-DNF}_n$ is efficiently PAC-learnable
2. Consistency problem for k -term DNF_n is NP-complete
3. 3-term $\text{DNF}_n \in 3\text{-CNF}_n$
4. 3-term DNF_n is efficiently PAC-learnable using 3-CNF
5. MONOM_n is efficiently PAC-learnable
6. Concept class of axis-aligned rectangles is eff. PAC-learn.
7. Concept class of symmetric boolean func is eff. PAC-learn.
8. $VC_{\text{dim}} = d+1$ for set of half-spaces in \mathbb{R}^d
9. $VC_{\text{dim}} = 3$ for circles in \mathbb{R}^2
10. $VC_{\text{dim}} = 2d+1$ for d -gons in \mathbb{R}^2

Infinite concept classes:

Family of subsets realized by C : $\Pi_C(S) = \{C \cap S: C \in C\}$

Number of realizable dichotomies: $\Pi_C(m) = \max_{S \subseteq \mathcal{X}, |S|=m} |\Pi_C(S)|$

Shattering: $S \subseteq \mathcal{X}$ is shattered by C , if $|\Pi_C(S)| = 2^{|S|}$

VC-Dimension d : Largest d , s.t. $\exists S \subseteq \mathcal{X}: |\Pi_C(S)| = |2^S| = d$

Error region: $\Delta(C) = \{h \in C: h \text{ is not consistent}\}$

ϵ -Net: $\Delta_\epsilon(C) = \{r \in \Delta(C): \Pr_{x \in D}[x \in r] \geq \epsilon\}$

Theorem (BEHW): C is PAC-learnable, if we draw m samples and find a consistent hypothesis with:

$$m \geq \frac{4}{\epsilon} \log_2 \frac{2}{\delta} + \frac{8 \cdot VC_{\text{dim}}(C)}{\epsilon} \log_2 \frac{13}{\epsilon} = O\left(\frac{1}{\epsilon} \left(\log \frac{1}{\delta} + VC_{\text{dim}}(C) \log \frac{1}{\epsilon}\right)\right)$$

Corollary:

1. Theorem holds
2. $VC_{\text{dim}}(C)$ is polynomial in n
3. Alg. runs polynomial in m and n

$\Rightarrow C$ is efficiently PAC-learnable

Theorem (PSSVC): For $d = VC_{\text{dim}}(C) > 0$:

1. $\Pi_C(m) = 2^m$ if $m \leq d$
2. $\Pi_C(m) \leq \left(\frac{me}{d}\right)^d$ else

Lemma: $C_1 \subseteq C_2 \Rightarrow VC_{\text{dim}}(C_1) \leq VC_{\text{dim}}(C_2)$

Lemma: For $\bar{C} = \{X \setminus c: c \in C\}$: $VC_{\text{dim}}(\bar{C}) = VC_{\text{dim}}(C)$

Radon's Theorem: Any set S with $d+2$ points in \mathbb{R}^d can be partitioned, s.t.:

$$\text{Conv}(S_1) \cap \text{Conv}(S_2) \neq \emptyset$$

Mangasarian Theorem: Given: 2 disjoint linearly separable subsets of \mathbb{R}^d

Then: Separating hyperplane can be found polynomial in d and $|S_1 \cup S_2|$

Optimization Theory

Constrained opt. problem: $\min_{\vec{w}} f(\vec{w})$ s.t. $g_i(\vec{w}) \leq 0$

Generalized Lagrangian: $L(\vec{w}, \vec{\lambda}) = f(\vec{w}) + \sum_{i=1}^m \lambda_i g_i(\vec{w})$

Dual opt. problem: $\max_{\vec{\lambda} \geq 0} \min_{\vec{w}} L(\vec{w}, \vec{\lambda})$

Decision Tree Learning

Entropy: $H(S) = -\sum_i P_i \log_2 \frac{1}{P_i}$

Lemma: $H(S) \leq \log_2(n)$

Cond. Entropy: $H(S|A) = -\sum_{j=1}^m P(A=U_j) \cdot H(S_j)$

Lemma: $H(S|A) \leq H(S)$

Information Gain: $H(S) - H(S|A)$

Gini Coefficient: $Gini(S) = 1 - (P_0^2 + P_1^2)$

Gini Split: $Gini_{\text{split}}(S, A) = \sum_{u \in \text{values}(A)} \frac{|S_u|}{|S|} Gini(S_u)$

Gini Gain: $Gini_{\text{gain}}(S, A) = Gini(S) - Gini_{\text{split}}(S, A)$

kernel Methods

kernel: $k(x, y) = \langle \phi(x), \phi(y) \rangle$

Gram/kernel Matrix $K_{X,K}$: $(K_{X,K})_{ij} = k(x_i, x_j)$

Finite Domain: k is kernel

$\Leftrightarrow K_{X,K}$ is positive semi-definite

Infinite Domain: (Mercer's Theorem)

k is kernel $\Leftrightarrow \forall \psi \in X: K_{\psi,k}$ is pos semi-definite

Cover's Theorem: Number of lin. separable dichotomies of n points in \mathbb{R}^d

$$C(n, d) = 2 \sum_{k=0}^{\lfloor \frac{n-1}{d} \rfloor} \binom{n-1}{k}$$

Prob. that a rand dich. is lin. sep.: $P(n, d) = \frac{1}{2^n} C(n, d)$

Phase transition: $n = 2(d+1)$

Kernels:

$$1. k(x, y) = f(x)f(y)$$

$$2. k(\vec{x}, \vec{y}) = \vec{x}^T A \vec{y}$$

$$3. k(x, y) = \alpha k_1(x, y) + \beta k_2(x, y)$$

$$4. k(x, y) = k_1(x, y)k_2(x, y)$$

$$5. k(x, y) = p(k_1(x, y))$$

$$6. k(x, y) = f(k_1(x, y))$$

$$7. k(x, y) = e^{k_1(x, y)}$$

$$8. k(\vec{x}, \vec{y}) = \exp\left(-\frac{\|\vec{x} - \vec{y}\|_2^2}{2\sigma^2}\right)$$

$$9. k(\vec{x}, \vec{y}) = (\vec{x}^T \vec{y} + c)^m$$

Ridge Regression:

least squares with L2-Tikhonov regularization

Optimization function:

$$\min_{\vec{w}} R_{\lambda}(\vec{w}, S) = \min_{\vec{w}} \|\vec{X}\vec{w} - \vec{y}\|^2 + \lambda \|\vec{w}\|^2$$

$$\text{Gradient: } \nabla_{\vec{w}} R_{\lambda}(\vec{w}, S) = 2\vec{X}^T \vec{X} \vec{w} - 2\vec{X}^T \vec{y} + 2\lambda \vec{w}$$

$$\text{Primal solution: } \vec{w} = (\vec{X}^T \vec{X} + \lambda \mathbf{I}_d)^{-1} \vec{X}^T \vec{y}$$

$$\text{Regression function: } f(\vec{x}) = \vec{x}^T \vec{w}$$

$$\text{Dual solution: } \vec{\alpha} = (\vec{X}\vec{X}^T + \lambda \mathbf{I}_n)^{-1} \vec{y}$$

$$\text{Regression function: } f(\vec{x}) = \sum_{i=1}^n \alpha_i \langle \vec{x}, \vec{x}_i \rangle$$

Lasso Regression

least squares with L2-Tikhonov regularization

Support Vector Machines

Hyperplane: $\langle \vec{w}, \vec{x} \rangle + b = 0$

$$\text{Margin: } \gamma = \frac{2}{\|\vec{w}\|}$$

Decision function: $\text{sign}(\langle \vec{w}, \vec{x} \rangle + b)$

γ -shattering: All dichotomies are realizable through a hyperplane with margin γ

Theorem: For any γ and any $S \subseteq \{\vec{x} \in \mathbb{R}^d : \|\vec{x}\| \leq R\}$

$$|S| \leq \min\left(\left(\frac{R}{\gamma}\right)^2, d\right) + 1$$

Support vectors: $\langle \vec{w}, \vec{x} \rangle + b = \pm 1, -1$

$$\alpha_i > 0$$

Hard Margin SVM

Primal form:

$$\max_{\vec{w}, b} \frac{2}{\|\vec{w}\|} \quad \text{with } y_i(\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1$$

$$\Leftrightarrow \min_{\vec{w}, b} \frac{1}{2} \|\vec{w}\|^2 \quad \text{with } -(y_i(\langle \vec{w}, \vec{x}_i \rangle + b) - 1) \leq 0$$

Dual form:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

$$\text{with } \sum_{i=1}^n \alpha_i y_i = 0 \text{ and } \alpha_i \geq 0$$

Maximum Margin Hyperplane:

$$f(\vec{x}) = \sum_{i=1}^n y_i \alpha_i \langle \vec{x}, \vec{x}_i \rangle + b \quad \text{with } b = -\frac{1}{2} \left(\max_{j: y_j = -1} \left(\sum_{i=1}^n y_i \alpha_i \langle \vec{x}_i, \vec{x}_j \rangle \right) + \min_{j: y_j = 1} \left(\sum_{i=1}^n y_i \alpha_i \langle \vec{x}_i, \vec{x}_j \rangle \right) \right)$$

$$\vec{w} = \sum_{i=1}^n y_i \alpha_i \vec{x}_i$$

Soft Margin SVM

Soft constraints: $y_i(\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - \xi_i$

$\xi_i = 0$: correct

$0 < \xi_i \leq 1$: within margin, but correct side

$\xi_i > 1$: wrong side

$$\text{Primal form: } \min_{\vec{w}, b, \xi} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\vec{w}\|^2$$

$$\text{with } y_i(\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

$$\text{Dual form: } \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

$$\text{with } 0 \leq \alpha_i \leq C \text{ and } \sum_{i=1}^n \alpha_i y_i = 0$$

Artificial Neural Networks

Perceptron rule:

$$\text{if } y_i(\vec{w}_k^T \vec{x}_i) \leq 0: \vec{w}_{k+1} = \vec{w}_k + y_i \vec{x}_i$$

Delta rule: $\vec{w} = \vec{w} + \Delta \vec{w}$, $\Delta \vec{w} = -\eta \nabla E(\vec{w})$

$$\text{Linear: } \frac{\partial E}{\partial w_i} = -\sum_{x \in \mathcal{D}} (y - \vec{w}^T \vec{x}) x_i$$

$$\text{Sigmoid: } \frac{\partial E}{\partial w_i} = -\sum_{x \in \mathcal{D}} (y - \sigma(\vec{z})) \sigma'(\vec{z}) (1 - \sigma(\vec{z})) x_i$$

Theorem: Block-Novikov

If: 1. $\forall i \in [n]: \|\vec{x}_i\| \leq R$ for some $R \in \mathbb{R}$

2. $\exists \epsilon \in \mathbb{R}^d, \gamma > 0: \forall i \in [n]: y_i(\vec{\epsilon}^T \vec{x}_i) \geq \gamma$

Then: #updates $k \leq \left(\frac{R}{\gamma}\right)^2$

Bayesian Learning

$$\begin{aligned} \text{MAP: } h_{\text{MAP}} &= \arg \max_{h \in \mathcal{H}} P(h|D) \\ &= \arg \max_{h \in \mathcal{H}} \frac{P(D|h)P(h)}{P(D)} \\ &= \arg \max_{h \in \mathcal{H}} P(D|h)P(h) \end{aligned}$$

$$\text{Likelihood: } L(h) = P(D|h) = \prod_{(x,y) \in D} P(y|x, h) P(x)$$

$$\text{Max. Likelihood } h: h_{\text{ML}} = \arg \max_{h \in \mathcal{H}} \sum_{(x,y) \in D} \log P(y|x, h)$$

$$\text{Bayes opt. class: } y_{\text{Bayes opt.}} = \arg \max_{y \in \mathcal{Y}} \sum_{h \in \mathcal{H}} P(y|h) P(h|x)$$

$$\text{Naive Bayes class: } y_{\text{NB}} = \arg \max_{y \in \mathcal{Y}} P(y) \prod_{i=1}^n P(a_i|y)$$

$$m\text{-estimate smoothing: } \hat{P}(a_i|y) = \frac{n_i + m p}{n + m}$$

$$\text{Laplace: } m = |\mathcal{A}|, p = \frac{1}{m}$$

Logistic Regression: (\Leftrightarrow Linear Regression)

$$\text{logit: } \log \frac{p}{1-p} \quad \text{prob: } \frac{1}{1+e^{-z}}$$

$$\text{Hypothesis: } \theta(\theta^T x') \quad \sigma = \frac{1}{1+e^{-z}}$$

Determine θ : Minimize:

$$L(\theta) = -\sum_{(x,y) \in D} y \log h_{\theta}(x) + (1-y) \log (1 - h_{\theta}(x))$$

$$\text{Partial derivatives: } \sum_{(x,y) \in D} \left(\frac{1}{1 + \exp(-\theta^T x')} - y \right) x_i$$

Comparing Learning Alg.

True error: $P_{x \in \mathcal{D}} [f(x) \neq h(x)]$

Sample error: $\frac{1}{|S|} \sum_{x \in S} \delta(f(x), h(x))$ $\delta = 1$, if $f(x) \neq h(x)$

Estimator: \hat{Y} estimates p of unknown dis.

Bias: $\mathbb{E}[\hat{Y}] - p$ Unbiased: Bias(\hat{Y}) = 0

$$\text{MSE: } \mathbb{E}[(\hat{Y} - p)^2]$$

$$\text{Variance: } \mathbb{E}[(\hat{Y} - \mathbb{E}[\hat{Y}])^2]$$

$$\text{Theorem: } \text{MSE}(\hat{Y}) = \text{Var}(\hat{Y}) + \text{Bias}(\hat{Y})^2$$

Central Limit Theorem:

$$\text{Sample mean: } \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$\text{For } n \rightarrow \infty: \bar{Y} \rightarrow \mu \quad \sigma \rightarrow \frac{\sigma}{\sqrt{n}}$$

True error estimator:

Estimator: Sample error

$$\text{error}_{S_i}(h) = \frac{1}{n}$$

Binom. Distribution:

$$P(r) = \binom{n}{r} \text{error}_D(h)^r (1 - \text{error}_D(h))^{n-r}$$

$$\text{Mean: } \mathbb{E}[X] = \sum_{i=0}^n P(i) = np$$

$$\text{Var.: } \text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = np(1-p)$$

$$\text{Std.: } \sigma_x = \sqrt{np(1-p)}$$

$$\text{Theorem: } \mathbb{E}[\text{error}_S(h)] - \text{error}_D(h) = 0$$

De Moivre-Laplace Theorem:

$$\mu = np \quad \sigma = \sqrt{np(1-p)}$$

Z-Scores: $z = \frac{\text{error}_S(h) - \text{error}_D(h)}{\sigma_{\text{error}_S(h)}} \sim \mathcal{N}(0, 1)$

7	50	68	90	95	98	99
2.17	1.1	1.28	1.64	1.96	2.33	2.58

$$\begin{aligned} 68\text{-}95\text{-}99.7 \text{ rule: } & 68.27\%: [\mu - \sigma, \mu + \sigma] \\ & 95.45\%: [\mu - 2\sigma, \mu + 2\sigma] \\ & 99.73\%: [\mu - 3\sigma, \mu + 3\sigma] \end{aligned}$$

Confidence interval for true error:

$$\text{error}_D(h) \in [\text{error}_S(h) - z \cdot \sigma, \text{error}_S(h) + z \cdot \sigma]$$

Difference between hypothesis:

$$\text{Estimator: } \hat{d} = \text{error}_{S_1}(h_1) - \text{error}_{S_2}(h_2)$$

$$\text{Approx.: } \mu = d$$

$$\sigma^2 \approx \frac{\text{error}_{S_1}(h_1)(1 - \text{error}_{S_1}(h_1))}{n_1} + \frac{\text{error}_{S_2}(h_2)(1 - \text{error}_{S_2}(h_2))}{n_2}$$

$$\text{Confidence interval: } \hat{d} \pm z_N \sigma$$

Comparing Learning Algorithms:

Estimate: $\mathbb{E}_{S \sim D} [\text{error}_D(L_A(S)) - \text{error}_D(L_B(S))]$

Paired t-Test:

Estimate: $\bar{S} = \mathbb{E}_{S \sim D} [\text{error}_D(L_A(S)) - \text{error}_D(L_B(S))]$

for S' with $|S'| = \frac{k-1}{k} |S|$

Confidence interval: $\bar{S} \pm t_{n,k-1} S_{\bar{S}}$

with: $S_{\bar{S}} = \sqrt{\frac{1}{k(k-1)} \sum_{i=1}^k (S_i - \bar{S})^2}$

$t_{n,k-1}$: t-distribution

Concept Learning:

Confusion matrix:

	predicted	
	+	-
actual +	TP	FN
actual -	FP	TN

Accuracy: $\frac{TP + TN}{TP + FN + FP + TN}$

Precision: $\frac{TP}{TP + FP}$

Recall: $\frac{TP}{TP + FN}$

F-score: $2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$

F-Beta score: $(1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}$

Hidden Markov Models

States: $S = \{S_1, S_2, \dots, S_N\}$

State at time t : $q_t \in S$

State Change: $\Pr(q_{t+1} = S_j | q_t = S_i, \dots, q_1 = S_1)$

Transition Probability Matrix A ($N \times N$)

$a_{ij} = \Pr(q_{t+1} = S_j | q_t = S_i) \geq 0$, s.t. $\sum_{j=1}^N a_{ij} = 1$

Initial Probabilities Π : $\Pi = (\pi_1, \dots, \pi_N)^T$

with $\pi_i = \Pr(q_1 = S_i)$

Emission Probability B ($N \times M$):

$b_j(w) = \Pr(O_t = w | q_t = S_j) \geq 0$

with: $\sum_{w=1}^M b_j(w) = 1$

Markov Assumption:

$$\Pr(q_{t+1} = S_j | q_t = S_i, \dots, q_1 = S_1) = \Pr(q_{t+1} = S_j | q_t = S_i)$$