

Supervised Learning

Phu Sakulwongtana

1 Introduction to Machine Learning Problem

Definition 1.1. (Machine Learning Problem) Define input space and output space $\mathcal{X} \subseteq \mathbb{R}^d$ and \mathcal{Y} , respectively. Given the training data points:

$$\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \subset \mathcal{X} \times \mathcal{Y}$$

The goal is to infer the function $f_S(\mathbf{x}_i) \approx y_i$, which we can use in the future data. There are 2 types of problems, when: $y \in \{-1, 1\}$, the problem is classification and if $y \in \mathbb{R}$, the problem is regression.

Definition 1.2. (Learning Algorithm) Given the training set $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \subset \mathcal{X} \times \mathcal{Y}$. The learning algorithm perform a mapping $S \mapsto f_S$, where the new input can be predicted as $f_S(\mathbf{x})$.

Definition 1.3. (Binary Classification) Given the training domain to be $\mathcal{X} = \mathbb{R}^2$ for $\mathbf{x} = (x_1, x_2)$ and $\mathcal{Y} = \{0, 1\}$, our predictor is defined as:

$$f(\mathbf{x}) = \begin{cases} 0 & \mathbf{w}^T \mathbf{x} + b > 0 \\ 1 & \mathbf{w}^T \mathbf{x} + b \leq 0 \end{cases}$$

Definition 1.4. (Mean Square Error) In most of the machine learning problem, we would like to find the predictor to minimize the following loss (for m size dataset):

$$\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x})^2$$

This is called mean-square error (MSE).

Lemma 1.1. Given the input and output dataset, which can be represented in matrix and vector notation:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

Given the predictor to be $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, then the mean-square error can be denoted as:

$$\mathcal{E}_{emp}(S, \mathbf{w}) = \frac{1}{m} (\mathbf{X}^T \mathbf{w} - \mathbf{y})^T (\mathbf{X}^T \mathbf{w} - \mathbf{y})$$

where the dataset is $S = (\mathbf{X}, \mathbf{w})$.

Proof. Consider the MSE to be, which we can consider the matrix multiplication:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 &= \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x})^2 \\ &= \frac{1}{m} \sum_{i=1}^m \left(y_i - \sum_{j=1}^n w_j x_{ij} \right)^2 \\ &= \frac{1}{m} (\mathbf{X}^T \mathbf{w} - \mathbf{y})^T (\mathbf{X}^T \mathbf{w} - \mathbf{y}) \end{aligned}$$

□

Proposition 1.1. *The solution to the mean-square error is given by:*

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Assuming that $\mathbf{X}^T \mathbf{X}$ is invertible.

Proof. Let's consider the derivative of $\nabla_{\mathbf{w}} \mathcal{E}_{\text{emp}}(S, \mathbf{w})$, which is given as:

$$\begin{aligned} \nabla_{\mathbf{w}} \left[(\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \right] &= \mathbf{0} \\ \Leftrightarrow \left(\sum_{i=1}^m \frac{\partial}{\partial w_1} \left(\sum_{j=1}^n x_{ij} w_j - y_i \right)^2, \dots, \sum_{i=1}^m \frac{\partial}{\partial w_n} \left(\sum_{j=1}^n x_{ij} w_j - y_i \right)^2 \right)^T &= \mathbf{0} \end{aligned}$$

Let's consider the derivative of each variable w_k

$$\begin{aligned} \frac{\partial \mathcal{E}_{\text{emp}}(S, \mathbf{w})}{\partial w_k} &= \frac{2}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i) \frac{\partial}{\partial w_k} \mathbf{w}^T \mathbf{x}_i \\ &= \frac{2}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ik} \end{aligned}$$

Let's consider the simpler case in 2 dimensions with $\mathbf{w} = (w_1, w_2)^T$, then setting this to zero gives us:

$$\sum_{i=1}^m (x_{ij} x_{i1} w_1 + x_{ik} x_{i2} w_2) = \sum_{i=1}^m x_{ik} y_i$$

for $k = 1, 2$. In vector notation, this is equivalent to $\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} = \sum_{i=1}^m \mathbf{x}_i y_i$ or it is equivalent to matrix notation is: $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$, taking the inverse gives us the required answer. □

Proposition 1.2. *Bias term for the predictor can be added i.e $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$.*

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{1} \\ \mathbf{1}^T \mathbf{X} & m \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{1}^T \mathbf{y} \end{bmatrix}$$

For dataset of size m , and $\mathbf{1}$ is the vector of elements 1.

Proof. This is equivalent to modify the dataset as $(\mathbf{x}^T, 1)$ with the same label y . Furthermore, the weight can be represented as (\mathbf{w}^T, b) . Now the linear equation (comes from the derivative) is:

$$\begin{aligned} (\mathbf{X}^T \mathbf{X}) \mathbf{w} &= \mathbf{X}^T \mathbf{1} b = \mathbf{X}^T \mathbf{y} \\ \mathbf{1}^T \mathbf{X} \mathbf{w} + m b &= \mathbf{1}^T \mathbf{y} \end{aligned}$$

This system of equation can be re-written as the matrix equation in the proposition, and so it is proven. □

Definition 1.5. (Nearest Neighbour) There are difference approach to training the predictor. We consider the set $N(\mathbf{x}; k)$ be the set of k -nearest (calculated using metrics) points to the point \mathbf{x} and its associated index set $I_{\mathbf{x}}$ i.e:

$$I_{\mathbf{x}} = \{i : \mathbf{x}_i \in N(\mathbf{x}; k)\}$$

The predictor function (for classification) is given by:

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } |\{y_i = 1 : i \in I_{\mathbf{x}}\}| > |\{y_i = 0 : i \in I_{\mathbf{x}}\}| \\ 0 & \text{otherwise} \end{cases}$$

On the other hand, the predictor for regression is defined by:

$$f(\mathbf{x}) = \frac{1}{k} \sum_{i \in I_{\mathbf{x}}} y_i$$

1.1 Bayes Estimator

Definition 1.6. (Expected Error) Assuming data is obtained by sampling iid from a fixed and unknown probability density $p(\mathbf{x}, y)$. The expected error of the predictor f is given by:

$$\mathcal{E}(f) = \mathbb{E} \left[(y - f(\mathbf{x}))^2 \right] = \iint (y - f(\mathbf{x}))^2 p(\mathbf{x}, y) \, d\mathbf{x} \, dy$$

Remark 1. The goal of our learning algorithm is to compute the optimal solution f^* as we have:

$$f^* = \arg \min_f \mathcal{E}(f)$$

However, to compute f^* , we have to know p . Please note that for the binary classification i.e where $\mathcal{Y} = \{0, 1\}$ and for given predictor f , the error $\mathcal{E}(f)$ is the average number of mistake of f .

Proposition 1.3. *The optimal solution f^* for regression problem $\mathcal{Y} = \mathbb{R}$, with square expected error. We can show that it is:*

$$f^*(\mathbf{x}) = \int_{\mathcal{Y}} y \, dp(y|\mathbf{x})$$

We assume that the joint distribution $p(y, \mathbf{x})$ can be decomposed as $p(y|\mathbf{x})p(\mathbf{x})$.

Proof. We have the following decomposition of the probability:

$$\mathcal{E}(f) = \int_{\mathcal{X}} \left\{ \int_{\mathcal{Y}} (y - f(\mathbf{x}))^2 \, dp(y|\mathbf{x}) \right\} \, dp(\mathbf{x})$$

We consider fixed $\mathbf{x} = \mathbf{x}'$ and given the following short-hand, we have $e = \mathcal{E}(f(\mathbf{x}'))$ and $z = f(\mathbf{x}')$, and so we have:

$$e \propto \int_{\mathcal{Y}} (y - z)^2 \, dp(y|\mathbf{x}')$$

The differentiation and setting this to zero giving us:

$$\begin{aligned} \frac{\partial e}{\partial z} &= -2 \int_{\mathcal{Y}} (y - z) \, dp(y|\mathbf{x}') \\ \iff 0 &= \int_{\mathcal{Y}} y \, dp(y|\mathbf{x}') - z \int_{\mathcal{Y}} dp(y|\mathbf{x}') \\ &= z - \int_{\mathcal{Y}} y \, dp(y|\mathbf{x}') \end{aligned}$$

This implies that $z = \int_{\mathcal{Y}} y \, dp(y|\mathbf{x}')$ and so the optimal predictor is equal to what we required. \square

1.2 Bias and Variance of Learning Algorithm

Remark 2. Assuming that there is a relationship between (\mathbf{x}, y) in the dataset, which is given by $y = F(\mathbf{x}) + \varepsilon$ where $\mathbb{E}[\varepsilon] = 0$ and finite variance. Then the optimal predictor can be shown to be:

$$f^*(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}] = F(\mathbf{x})$$

Remark 3. We want to consider the **expected error by an arbitrary learner $A_{\mathcal{S}}(\mathbf{x})$** . The expected error is:

$$\mathcal{E}(A_{\mathcal{S}}(\mathbf{x}')) = \mathbb{E}[(y' - A_{\mathcal{S}}(\mathbf{x}))^2]$$

where y' is sample from the marginal $p(y|\mathbf{x}')$.

Lemma 1.2. *We can show that:*

$$\mathbb{E}[(Z - \mathbb{E}[X])^2] = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$$

Proof. We have the following:

$$\begin{aligned}\mathbb{E}[(Z - \mathbb{E}[Z])^2] &= \mathbb{E}[Z^2 - 2Z + \mathbb{E}[Z]^2] \\ &= \mathbb{E}[Z^2] - 2\mathbb{E}[Z]^2 + \mathbb{E}[Z]^2 \\ &= \mathbb{E}[Z^2] - \mathbb{E}[Z]^2\end{aligned}$$

□

Proposition 1.4. (Decomposing) The square error $\mathcal{E}(A(\mathbf{x}'))$ can be decomposed to:

$$\mathbb{E}[(y - f^*(\mathbf{x}'))^2] + (f^*(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2 + \mathbb{E}[(A_S(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2]$$

Proof. We can decomposed the error of the learner $\mathcal{E}(A_S(\mathbf{x}'))$ as we have:

$$\begin{aligned}\mathbb{E}[(y' - A_S(\mathbf{x}'))^2] &= \mathbb{E}[(y')^2 - 2y'A_S(\mathbf{x}') + A_S(\mathbf{x}')^2] \\ &= \mathbb{E}[(y' - f^*(\mathbf{x}'))^2] + f^{*2}(\mathbf{x}') - 2f^*(\mathbf{x}')\mathbb{E}[A_S(\mathbf{x}')] + \mathbb{E}[(A_S(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2] + \mathbb{E}[A_S(\mathbf{x}')]^2 \\ &= \mathbb{E}[(y - f^*(\mathbf{x}'))^2] + (f^*(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2 + \mathbb{E}[(A_S(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2]\end{aligned}$$

Let's show that the second equality is actually equal to the first equation:

$$\begin{aligned}&\mathbb{E}[(y' - \mathbb{E}[y'|\mathbf{x}'])^2] + \mathbb{E}[y'|\mathbf{x}']^2 - 2\mathbb{E}[y'|\mathbf{x}']\mathbb{E}[A_S(\mathbf{x}')] + \mathbb{E}[(A_S(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2] + \mathbb{E}[A_S(\mathbf{x}')]^2 \\ &= \mathbb{E}[(y')^2] - \mathbb{E}[y'|\mathbf{x}']^2 + \mathbb{E}[y'|\mathbf{x}']^2 - 2\mathbb{E}[y'|\mathbf{x}']\mathbb{E}[A_S(\mathbf{x}')] + \mathbb{E}[A_S(\mathbf{x}')^2] - \mathbb{E}[A_S(\mathbf{x}')]^2 + \mathbb{E}[A_S(\mathbf{x}')]^2 \\ &= \mathbb{E}[(y')^2] - 2\mathbb{E}[y'|\mathbf{x}']\mathbb{E}[A_S(\mathbf{x}')] + \mathbb{E}[A_S(\mathbf{x}')^2] \\ &= \mathbb{E}[(y' - A_S(\mathbf{x}'))^2]\end{aligned}$$

As required. □

Remark 4. (Bias and Variance Tradeoff) We can see that each term in the decomposition has the following contribution:

$$\underbrace{\mathbb{E}[(y - f^*(\mathbf{x}'))^2]}_{\text{Bayes' Error}} + \underbrace{(f^*(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(A_S(\mathbf{x}') - \mathbb{E}[A_S(\mathbf{x}')])^2]}_{\text{Variance}}$$

The **bias error** describes the discrepancy between the algorithm and truth value. The **Bayes error** is the **irreducible noise**. Finally, **variance capture the variance of the algorithm between training set**. We can have the additional observation:

- Bias and Variance tends to trade-off against one another.
- Many parameter allows better flexibility to fit the data, which lower the bias. However, it also gives rise to high-variance, and vice versa.
- This composition holds for square loss function.

Definition 1.7. (Bayes Estimator for Classification) For C-class classification (Bayes classifier), it is given by:

$$f^*(\mathbf{x}) = \arg \max_{c \in [C]} p(y = c|\mathbf{x})$$

where the loss is 0 if we predict correctly and 1 otherwise. Furthermore, the Bayes optimal error rate is:

$$\int \left(1 - p(y = f^*(\mathbf{x})|\mathbf{x})\right) dp(\mathbf{x})$$

Lemma 1.3. For Z being a random variable with values $[0, 1]$ and let $\mathbb{E}[Z] = \mu$ for any $a \in (0, 1)$ we have:

$$\mathbb{P}(Z > 1 - a) > \frac{\mu - (1 - a)}{a}$$

Proof. Recall the Markov's inequality:

$$\begin{aligned}\mathbb{P}(Z \geq a) &\leq \frac{\mathbb{E}[Z]}{a} \\ \implies 1 - \mathbb{P}(Z \geq a) &\geq 1 - \frac{\mathbb{E}[Z]}{a}\end{aligned}$$

We consider the following inequality, where we consider:

$$\begin{aligned}\mathbb{P}(1 - Z < a) &\geq 1 - \frac{\mathbb{E}[1 - Z]}{a} \\ &= 1 - \frac{1 - \mu}{a} = \frac{a - 1 + \mu}{a} = \frac{\mu - (1 - a)}{a}\end{aligned}$$

□

Theorem 1.1. (No Free-Lunch) Let A be any learning algorithm for binary classifier (where $\mathcal{Y} = \{-1, 1\}$) over domain \mathcal{X} . Let $m < |\mathcal{X}|/2$ being a training size. We define the loss of the function f to be:

$$\mathcal{E}_p(f) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{I}[f(\mathbf{x}) \neq y] \, dp(\mathbf{x}, y)$$

Then there exists a distribution p such that:

- There exists a function $f : \mathcal{X} \rightarrow \{0, 1\}$ with $\mathcal{E}_p(f) = 0$
- For dataset $\mathcal{S} \sim p^m$, we have

$$\mathbb{P}_{\mathcal{S} \sim p^m} [\mathcal{E}_p(A(\mathcal{S})) > 1/8] \geq 1/7$$

Proof. This prove is abit more involved. Let's start with proving the first point (For now we assume the discrete distribution and finite value of \mathcal{X}). Let $C \subset \mathcal{X}$, where $|C| = 2m$. Denote \mathcal{Y}^C being the set of all possible function $f : C \rightarrow \mathcal{Y}$ i.e $\{f_1, \dots, f_T\}$ where $T = 2^{2m}$. We can construct the distribution function p_i such that:

$$p_i(\{\mathbf{x}, y\}) = \begin{cases} 1/(2m) & \text{if } y = f_i(\mathbf{x}) \\ 0 & \text{otherwise} \end{cases}$$

for all \mathbf{x}, y . Let's consider $\mathcal{E}_{p_i}(f_i)$, which is:

$$\begin{aligned}\mathcal{E}_{p_i}(f_i) &= \sum_{\mathbf{x} \in \mathcal{X}} \sum_{y \in \{-1, 1\}} \mathbb{I}[f_i(\mathbf{x}) \neq y] p_i(\{\mathbf{x}, y\}) \\ &= \sum_{\mathbf{x} \in C} \mathbb{I}[f_i(\mathbf{x}) \neq f_i(\mathbf{x})] = 0\end{aligned}$$

And so the first point is proven. Now, consider all possible combination of data points of size m in C i.e $C^m = \{S_1, \dots, S_k\}$ where $k = (2n)^n$. We construct the dataset from f_i as $S_j^i = \{(\mathbf{x}, f_i(\mathbf{x})) : \mathbf{x} \in S_j\}$. Now consider the expected error of an algorithm under correction function f_i i.e

$$\begin{aligned}\mathbb{E}_{\mathcal{S} \sim p_i^m} [\mathcal{E}_{p_i}(A(\mathcal{S}))] &= \sum_{j=1}^k p_i(S_j^i) \mathcal{E}_{p_i}(A(S_j^i)) \\ &= \sum_{j=1}^k \frac{1}{(2n)^n} \mathcal{E}_{p_i}(A(S_j^i)) = \frac{1}{k} \sum_{j=1}^k \mathcal{E}_{p_i}(A(S_j^i))\end{aligned}$$

Note that for scalar $\alpha_1, \dots, \alpha_m$ we have $\max_l \alpha_l \geq 1/m \sum_{i=1}^m \alpha_i$ and $\min_l \alpha_l \leq 1/m \sum_{i=1}^m \alpha_i$. Consider the value of the function f_i that maximizes the error of the learner (when everything is under f_i):

$$\begin{aligned} \max_{i \in [T]} \mathbb{E}_{S \sim p_i^m} [\mathcal{E}_{p_i}(A(S))] &\geq \frac{1}{T} \sum_{i=1}^T \frac{1}{k} \sum_{j=1}^k \mathcal{E}_{p_i}(A(S_j^i)) \\ &= \frac{1}{k} \sum_{j=1}^k \frac{1}{T} \sum_{i=1}^T \mathcal{E}_{p_i}(A(S_j^i)) \geq \min_{j \in [k]} \frac{1}{T} \sum_{i=1}^T \mathcal{E}_{p_i}(A(S_j^i)) \end{aligned}$$

Now, denote a set $S'_j = \{\mathbf{v}_1, \dots, \mathbf{v}_p\} \subset C$ such that its element doesn't belong in S_j for $j = 1, \dots, k$, consider average expected risk:

$$\begin{aligned} \frac{1}{T} \sum_{i=1}^T \mathcal{E}_{p_i}(A(S_j^i)) &= \frac{1}{T} \sum_{i=1}^T \sum_{\mathbf{x} \in \mathcal{X}} \sum_{y \in \{-1, 1\}} \mathbb{I}[A(S_j^i)(\mathbf{x}) \neq y] \rho_i(\{\mathbf{x}, y\}) \\ &\geq \frac{1}{T} \sum_{i=1}^T \sum_{\mathbf{v} \in S'_j} \sum_{y \in \{-1, 1\}} \mathbb{I}[A(S_j^i)(\mathbf{v}) \neq y] \rho_i(\{\mathbf{v}, y\}) \\ &= \frac{1}{T} \sum_{i=1}^T \sum_{\mathbf{v} \in S'_j} \mathbb{I}[A(S_j^i)(\mathbf{v}) \neq f_i(\mathbf{v})] \rho_i(\{\mathbf{v}, f_i(\mathbf{v})\}) \\ &= \frac{1}{T} \sum_{i=1}^T \sum_{\mathbf{v} \in S'_j} \frac{1}{2m} \mathbb{I}[A(S_j^i)(\mathbf{v}) \neq f_i(\mathbf{v})] \\ &\geq \frac{1}{T} \sum_{i=1}^T \sum_{\mathbf{v} \in S'_j} \frac{1}{2p} \mathbb{I}[A(S_j^i)(\mathbf{v}) \neq f_i(\mathbf{v})] \\ &= \frac{1}{2} \frac{1}{p} \sum_{\mathbf{v} \in S'_j} \frac{1}{T} \sum_{i=1}^T \mathbb{I}[A(S_j^i)(\mathbf{v}) \neq f_i(\mathbf{v})] \\ &= \frac{1}{2} \min_{r \in [p]} \frac{1}{T} \sum_{i=1}^T \mathbb{I}[A(S_j^i)(\mathbf{v}_r) \neq f_i(\mathbf{v}_r)] \end{aligned}$$

Note that $p \geq m$ because the dataset doesn't have to be unique. Before further analysis, for \mathcal{Y}^C , we can partition into $T/2$ pairs $(f_i, f_{i'})$ such that $f_i(\mathbf{x}) \neq f_{i'}(\mathbf{x})$ iff $\mathbf{x} = \mathbf{v}_r$ for $r \in [p]$, by setting $f_{i'}(\mathbf{v}_r) = \neg f_i(\mathbf{v}_r)$ where

$$\neg a = \begin{cases} 1 & \text{if } a = -1 \\ -1 & \text{if } a = 1 \end{cases}$$

Please note that $S_j^i = S_j^{i'}$ because the effect of $f_{i'}(x) \neq f_i(x)$ iff $x \notin S_j^i$. Thus, we can see that:

$$\mathbb{I}[A(S_j^i)(\mathbf{v}_r) \neq f_i(\mathbf{v}_r)] + \mathbb{I}[A(S_j^{i'})(\mathbf{v}_r) \neq f_{i'}(\mathbf{v}_r)] = 1$$

Let's consider the value inside, by the partition of list of all functions, we have:

$$\begin{aligned} \frac{1}{T} \sum_{i=1}^T \mathbb{I}[A(S_j^i)(\mathbf{v}_r) \neq f_i(\mathbf{v}_r)] &= \frac{1}{T} \sum_{(i, i')} \mathbb{I}[A(S_j^i)(\mathbf{v}_r) \neq f_i(\mathbf{v}_r)] + \mathbb{I}[A(S_j^{i'})(\mathbf{v}_r) \neq f_{i'}(\mathbf{v}_r)] \\ &= \frac{1}{T} \sum_{(i, i')} 1 = \frac{1}{T} \frac{T}{2} = \frac{1}{2} \end{aligned}$$

This implies that:

$$\max_{i \in [T]} \mathbb{E}_{S \sim p_i^m} [\mathcal{E}_{p_i}(A(S))] \geq \min_{j \in [k]} \frac{1}{T} \sum_{i=1}^T \mathcal{E}_{p_i}(A(S_j^i)) \geq \frac{1}{4}$$

This means that for all algorithm A' getting a dataset S of size m , there is a function f and a distribution p over $\mathcal{X} \times \{0, 1\}$ such that:

$$\mathbb{E}_{S \sim p^m}[\mathcal{E}_p(A(S))] \geq \frac{1}{4}$$

and so, using the probabilistic inequality, we have:

$$\begin{aligned} \mathbb{P}\left[\mathcal{E}_p(A(S)) \geq \frac{1}{8}\right] &= \mathbb{P}\left[\mathcal{E}_p(A(S)) \geq 1 - \frac{7}{8}\right] \geq \frac{\mathbb{E}_{S \sim p^m}[8\mathcal{E}_p(A(S))](1 - 7/8)}{7/8} \\ &\geq \frac{2 - 1}{7} = \frac{1}{7} \end{aligned}$$

Thus complete the proof. \square

Theorem 1.2. *As the number of sample goes to infinity, the error rate is no more than twice of the Bayes error rate for the k -nearest neighbour. Please note that the k -nearest neighbour attempts to approximate:*

$$p(y = c|\mathbf{x}) \approx \frac{|\{i : y_i = c, i \in I_{\mathbf{x}}\}|}{k}$$

We consider the points that are near the evaluation points and find the class of the neighbours that has the highest frequency.

Proof. (Sketch) We will shorten the notation as $p(c|\mathbf{x}) = p(y = c|\mathbf{x})$. The expected Bayes classifier (at \mathbf{x}) is:

$$1 - \max_{c \in [C]} p(c|\mathbf{x})$$

The expected error rate of 1-NN at \mathbf{x} is given by:

$$\sum_{c=1}^C p_{\text{nn}}(c|\mathbf{x})[1 - p(c|\mathbf{x})]$$

As the number of sequence goes to infinity $m \rightarrow \infty$, we have $p(c|\mathbf{x}) \approx p_{\text{nn}}(c|\mathbf{x})$. Now, we will show that:

$$\sum_{c=1}^C p(c|\mathbf{x})[1 - p(c|\mathbf{x})] \leq 2 \left[1 - \max_{c \in [C]} p(c|\mathbf{x}) \right]$$

Let $c^* = \arg \max_{c \in [C]} p(c|\mathbf{x})$ and $p^* = p(c^*|\mathbf{x})$ observe that $c \in [C]$:

$$\begin{aligned} \sum_{c=1}^C p(c|\mathbf{x})[1 - p(c|\mathbf{x})] &= p^*(1 - p^*) + \sum_{c \in [C] \setminus c^*} p(c|\mathbf{x})[1 - p(c|\mathbf{x})] \\ &\leq (C - 1) \frac{1 - p^*}{C - 1} \left[1 - \frac{1 - p^*}{C - 1} \right] + p^*(1 - p^*) \\ &= (1 - p^*) \left[1 - \frac{1 - p^*}{C - 1} + p^* \right] \\ &\leq (1 - p^*) [1 + p^*] \leq 2(1 - p^*) \end{aligned}$$

The second inequality comes from the fact that sum is maximized when all $p(c|\mathbf{x})$ have the same value. And the last inequality comes from the fact that $p, p^* < 1$. Thus complete the proof. \square

Remark 5. One can show that for $k = k(m)$ where m is the size of the dataset, one can show that

$$\mathcal{E}(k - \text{NN}) \rightarrow \mathcal{E}(f^*)$$

note that k depends on the data size, as $m \rightarrow \infty$ with the condition that $k(m) \rightarrow \infty$ and $k(m)/m \rightarrow 0$.

Remark 6. (Curse of Dimensionality) The rate of convergence depends exponentially on the input dimension. This problem occurs throughout the ML algorithms. The intuition is that the volume increases exponentially with the dimension implies that the number of data required to cover the space to perform estimate also increase exponentially: The ratio between volume unit d -dim ball centered at the origin and $1/2$ -unit ball at the origin is $(1/2)^d$.

Definition 1.8. (Empirical Risk) We are given only the sample from probability $p(\mathbf{x}, y)$. The natural approach is to approximate the expected error using empirical error:

$$\mathcal{E}_{\text{emp}}(\mathcal{S}, f) = \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{x}, y) \in \mathcal{S}} (y_i f(\mathbf{x}_i))^2$$

Definition 1.9. (Empirical Risk Minimization (with regularizer)) If we consider all possible function, we can always find the function with 0 empirical error (remember) this is known as overfitting. To solve this we have to restrict the function space to be \mathcal{H} called hypothesis space, which ERM is defined as:

$$f_{\mathcal{S}} = \arg \min_{f \in \mathcal{H}} \mathcal{E}_{\text{emp}}(\mathcal{S}, f)$$

Remark 7. (Example of Hypothesis Space) We can consider the following increasing “complexity” as for the regression in 1D as we have:

$$\mathcal{H}_n = \left\{ f(x) = \sum_{l=1}^n a_l x^l + b : a_1, \dots, a_n, b \in \mathbb{R} \right\}$$

Choosing the correct model requires a cross-validation. Unless the prior knowledge is available on f^* , we can't expect $f^* \in \mathcal{H}$, while we can't allow too large \mathcal{H} as it leads to the overfitting.

2 Kernel and Regression

2.1 Introduction

Definition 2.1. (Convex Set) A set \mathcal{X} is convex if $\mathbf{p}, \mathbf{q} \in \mathcal{X}$ is convex if $\alpha \mathbf{p} + (1 - \alpha) \mathbf{q} \in \mathcal{X}$

Definition 2.2. (Convex Function) A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is convex iff for all $\mathbf{p}, \mathbf{q} \in \mathcal{X}$ in convex set and $\alpha \in (0, 1)$ as we have:

$$f(\alpha \mathbf{p} + (1 - \alpha) \mathbf{q}) \leq \alpha f(\mathbf{p}) + (1 - \alpha) f(\mathbf{q})$$

A function f is concave if $-f$ is convex. A function is *strictly convex* if we replace \leq with $<$.

Remark 8. (Various Comments on Convex Function) We have the following results on the convex function, as we have:

- If f and g are convex, then $f + g$ is convex.
- If f is convex and g is affine (linear + constant) then $f(g(\cdot))$ is convex.
- Suppose \mathbf{M} is symmetric matrix, then \mathbf{M} is positive semi-definite matrix iff $f(\mathbf{x}) = \mathbf{x}^T \mathbf{M} \mathbf{x}$ is convex.
- Level set $\{\mathbf{x} : f(\mathbf{x}) = c\}$ where $c \in \mathbb{R}$ of convex function f is convex.
- For $f : (a, b) \rightarrow \mathbb{R}$ if $f'' \geq 0$ then f is convex.
- For $f : \mathcal{X} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ if $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$ for all $\mathbf{x} \in \mathcal{X}$, then f is convex.

2.2 Ridge Regression

Definition 2.3. (Ridge Regression Problem) Given a function $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ with a dataset:

$$\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \subset \mathbb{R}^n \times \mathbb{R}$$

Assuming the dataset is generated by the unknown function g i.e $(\mathbf{x}, g(\mathbf{x}))$. Then suppose that the vector \mathbf{x}_i are linearly independent with $m = n$, then there is a unique solution, whose parameter \mathbf{w} solves:

$$\mathbf{X}\mathbf{w} = \mathbf{y}$$

where $\mathbf{y} = (y_1, \dots, y_m)^T$ and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]^T \in \mathbb{R}^{m \times n}$.

Definition 2.4. (Well-Posed) The solution/problem is called well-posed if: the solution exists, unique and depends continuously on the data. The regularized theory allows general framework to solve ill-posed problem (we can choose the term to penalize complex function).

Definition 2.5. (Regularized Empirical Error) We minimize the following regularized empirical error, which is given by:

$$\begin{aligned} \mathcal{E}_{\text{emp}, \lambda}(\mathbf{w}) &= \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{i=1}^n w_i^2 \\ &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \end{aligned}$$

We can see that the parameter $\lambda > 0$ defines the trade-off between error and the norm of vector \mathbf{w} (which restricts the complexity of the model).

Proposition 2.1. Solving the regularized empirical error by setting its gradient to $\mathbf{0}$, gives us:

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_n)^{-1} \mathbf{X}^T \mathbf{y}$$

Furthermore, we can show that the weight $\mathbf{w} = \sum_{i=1}^m \alpha_i \mathbf{x}_i$ and the solution can be written as:

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i \mathbf{x}_i^T \mathbf{x}_i$$

where $\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_m)^{-1} \mathbf{y}$. This is called dual form, while $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ is called primal form.

Proof. Starting with the derivative, we have:

$$\nabla \mathcal{E}_{\text{emp}, \lambda}(\mathbf{w}) = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + 2\lambda \mathbf{w} = \mathbf{0}$$

which implies the weight of the first form i.e $\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_n)^{-1} \mathbf{X}^T \mathbf{y}$. Now, we can also see that:

$$\mathbf{w} = \frac{\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{\lambda}$$

Assume the the dual form of the weight $\mathbf{w} = \sum_{i=1}^m \alpha_i \mathbf{x}_i$ as we have:

$$\alpha_i = \frac{y_i - \mathbf{w}^T \mathbf{x}_i}{\lambda} = \frac{y_i - (\sum_{i=1}^m \alpha_i \mathbf{x}_i)^T \mathbf{x}_i}{\lambda}$$

Now solving for the value of y_i , which we have:

$$\begin{aligned} y &= \left(\sum_{j=1}^m \alpha_j \mathbf{x}_j \right)^T \mathbf{x}_i + \lambda \alpha_i \\ &= \sum_{j=1}^m (\alpha \mathbf{x}_j^T \mathbf{x}_j + \lambda \alpha_j \delta_{ij}) = \sum_{j=1}^m (\mathbf{x}_j^T \mathbf{x}_j + \lambda \delta_{ij}) \alpha \end{aligned}$$

and so we have $(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_m) \boldsymbol{\alpha} = \mathbf{y}$

□

Remark 9. (Advantage of Dual Form) The dual form allow us to gain a computational advantage for both training and testing time:

- *Training Time:* Solving \mathbf{w} in the primal function requires $\mathcal{O}(mn^2 + n^3)$ operations while solving for dual form $\mathcal{O}(nm^2 + m^3)$ if $m \ll n$ then it is more efficient than primal.
- *Testing Time:* Computing $f(\mathbf{x})$ in test vector \mathbf{x} in the primal form requires $\mathcal{O}(n)$ operations but the dual form requires $\mathcal{O}(nm)$ operations.

2.3 Basis/Kernel Functions

Definition 2.6. (Basis/Feature Function) We have the function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^N$ as we have:

$$\phi(\mathbf{x}) = \left(\phi_1(\mathbf{x}), \dots, \phi_N(\mathbf{x}) \right)^T$$

for $\mathbf{x} \in \mathbb{R}^n$, where we call ϕ_1, \dots, ϕ_N are called basis function and $\phi(\mathbf{x})$ is called feature vector, and feature space is defined by: $\{\phi(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}$

Remark 10. We can use the feature map of the data instead of real data $\phi(\mathbf{x})$. This gives us the many advantages, for example:

- The map: $\phi(\mathbf{x}) = (\mathbf{x}, 1)^T$ allow us to have the bias terms.
- The map: $\phi(\mathbf{x}) = (\mathbf{x}_1 x_2)^T$ allow us to consider the interaction between inputs (individual elements).

We can also consider the second order correlation if $\mathbf{x} \in \mathbb{R}^n$ as:

$$\phi(\mathbf{x}) = (x_1 x_1, x_1 x_2, \dots, x_1 x_n, x_2 x_2, x_2 x_3, \dots, x_2 x_n, \dots, x_n x_n)^T$$

now the feature vector has the size of $(n^2 + n)/2$. However, if we consider the inner product, we will have:

$$\begin{aligned} \langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle &= (x_1 x_1, x_1 x_2, \dots, x_n x_n)^T (t_1 t_1, t_1 t_2, \dots, t_n t_n) \\ &= (x_1 t_1 + \dots + x_n t_n)(x_1 t_1 + \dots + x_n t_n) \\ &= (\mathbf{x}^T \mathbf{t})^T \end{aligned}$$

Note that $\mathcal{O}(n)$ but the native computation will take $\mathcal{O}(n^2)$. This leads to decrease the computation complexity (please see the dual form too).

Definition 2.7. (Kernel Function) Given a feature map ϕ , we define the associated kernel function $k : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ as we have:

$$k(\mathbf{x}, \mathbf{t}) = \langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle$$

Please note that the computing $k(\mathbf{x}, \mathbf{t})$, which it doesn't depends on computing $\phi(\mathbf{x})$.

Remark 11. (Feature Map not Unique) The feature map isn't unique. Consider the ϕ that is associated with kernel k , and so $\hat{\phi} = U\phi$ where $U \in \mathbb{R}^{N \times N}$. The feature can be difference in values and dimension but gives rise to the same kernel:

$$(U\phi)^T (U\phi) = \phi^T \phi$$

Theorem 2.1. (Representer) Consider the loss to be:

$$\mathcal{E}_{emp, \lambda}(\mathbf{w}) = \sum_{i=1}^m V(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle) + \lambda \langle \mathbf{w}, \mathbf{w} \rangle$$

where $V : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is a loss function. If V is differentiable with respect to its second argument and \mathbf{w} is a minimizer of \mathcal{E}_λ , then \mathbf{w} has the form of:

$$\mathbf{w} = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i) \implies f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle = \sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

Proof. The proof is similar to the dual form. Setting the derivative of \mathcal{E}_λ with respect to zero and we have:

$$\sum_{i=1}^m V'(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle) \phi(\mathbf{x}_i) + 2\lambda \mathbf{w} = 0$$

Compared to $\mathbf{w} = \sum_{i=1}^m \alpha_i \phi(\mathbf{x}_i)$, we can see that:

$$\alpha_i = \frac{1}{2\lambda} V'(y_i, \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle)$$

From the definition of \mathbf{w} , we can see that:

$$\alpha_i = \frac{1}{2\lambda} V' \left(y_i, \sum_{j=1}^m k(\mathbf{x}_i, \mathbf{x}_j) \alpha_j \right)$$

for $i = 1, \dots, m$. Finding $\boldsymbol{\alpha}$ is done by solving the following optimization problem

$$\arg \max_{\boldsymbol{\alpha}} \sum_{i=1}^m V(y_i, (\mathbf{K}\boldsymbol{\alpha})_i) + \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha}$$

□

Definition 2.8. (Positive Semi-Definite Kernel) The kernel $k : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is positive semi-definite if it is symmetric and given the set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, the matrix:

$$\begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

is positive semi-definite.

Theorem 2.2. Kernel k is positive definite iff:

$$k(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ for some feature map $\phi : \mathbb{R}^n \rightarrow \mathcal{W}$ for Hilbert space \mathcal{W}

Proof. We will consider only one direction. If $k(\mathbf{x}, \mathbf{t}) = \langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle$, then we have:

$$\sum_{i=1}^n \sum_{j=1}^m c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) = \left\langle \sum_{i=1}^n c_i \phi(\mathbf{x}_i), \sum_{j=1}^m c_j \phi(\mathbf{x}_j) \right\rangle = \left\| \sum_{i=1}^n c_i \phi(\mathbf{x}_i) \right\|^2 \geq 0$$

□

Definition 2.9. (Polynomial Kernel) If $p : \mathbb{R} \rightarrow \mathbb{R}$ is a polynomial with non-negative coefficient then $k(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}^T \mathbf{z})$ where $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$ and k positive semi-definite kernel.

Proposition 2.2. If \mathbf{A} is an $n \times n$ positive semi-definite matrix, the function $k : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ defined by:

$$k(\mathbf{x}, \mathbf{t}) = \mathbf{x}^T \mathbf{A} \mathbf{t}$$

is a generalized linear kernel and it is a positive semi-definite kernel.

Proof. Since \mathbf{A} is positive semi-definite, we can write \mathbf{A} in the form of $\mathbf{A} = \mathbf{R}\mathbf{R}^T$ for some $\mathbf{R} \in \mathbb{R}^{n \times n}$. Thus, k is represented by a feature map $\phi(\mathbf{x}) = \mathbf{R}^T \mathbf{x}$. As we can see that:

$$\begin{aligned} \sum_{ij} c_i c_j \mathbf{x}_i^T \mathbf{A} \mathbf{x}_j &= \sum_{ij} c_i c_j (\mathbf{R}^T \mathbf{x}_i)^T (\mathbf{R}^T \mathbf{x}_j) \\ &= \sum_i c_i [\mathbf{R}^T \mathbf{x}_i]^T \left[\sum_j c_j (\mathbf{R}^T \mathbf{x}_j) \right] = \left\| \sum_i c_i \mathbf{R}^T \mathbf{x}_i \right\|^2 \geq 0 \end{aligned}$$

□

Proposition 2.3. If $k : \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}$ is a positive semi-definite kernel and $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^N$.

$$\tilde{k}(\mathbf{x}, \mathbf{t}) = k(\phi(\mathbf{x}), \phi(\mathbf{t}))$$

The kernel \tilde{k} defined to be $\tilde{k} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a positive definite kernel.

Proposition 2.4. Given a positive semi-definite kernels k_1 and k_2 , ak_1 is a positive semi-definite kernel if $a > 0$ and $k_1 + k_2$ is also a positive definite kernel.

Proposition 2.5. We consider the following combination of kernel k_1 and k_2 are given as:

$$k(\mathbf{x}, \mathbf{t}) = k_1(\mathbf{x}, \mathbf{t})k_2(\mathbf{x}, \mathbf{t})$$

where $\mathbf{x}, \mathbf{t} \in \mathbb{R}^d$ is a kernel.

Proof. For the product of kernel, we have:

- We want to show that for positive semi-definite \mathbf{A} and \mathbf{B} where $\mathbf{C} = \mathbf{A} \odot \mathbf{B}$ is a positive semi-definite.
- Since \mathbf{A} and \mathbf{B} are positive semi-definite, where it can be factorized as $\mathbf{A} = \mathbf{U}\mathbf{U}^T$ and $\mathbf{B} = \mathbf{V}\mathbf{V}^T$ for $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times n}$ as we have:

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n z_i z_j C_{ij} &= \sum_{i=1}^n \sum_{j=1}^n z_i z_j \left(\sum_{r=1}^n U_{ir} U_{jr} \right) \left(\sum_{s=1}^n V_{is} V_{js} \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n \sum_{r=1}^n \sum_{s=1}^n z_i z_j U_{ir} U_{jr} V_{is} V_{js} \\ &= \sum_{r=1}^n \sum_{s=1}^n \sum_{i=1}^n \sum_{j=1}^n z_i z_j U_{ir} U_{jr} V_{is} V_{js} \\ &= \sum_{r=1}^n \sum_{s=1}^n \sum_{i=1}^n z_i U_{ir} V_{is} \sum_{j=1}^n z_j U_{jr} V_{js} = \sum_{r=1}^n \sum_{s=1}^n \left(\sum_{i=1}^n z_i U_{ir} V_{is} \right)^2 \geq 0 \end{aligned}$$

Thus complete the proof. This proves the polynomial kernel is positive definite kernel.

□

Remark 12. (Several Kernels) We have the following positive definite kernel, where we have $a \geq 0$:

- $k(\mathbf{x}, \mathbf{t}) = (\mathbf{x}^T \mathbf{t})^r$
- $k(\mathbf{x}, \mathbf{t}) = (a + \mathbf{x}^T \mathbf{t})^r$
- $k(\mathbf{x}, \mathbf{t}) = \sum_{i=1}^d (a^i / i!) (\mathbf{x}^T \mathbf{t})^r$

- Gaussian Kernel: $k(\mathbf{x}, \mathbf{t}) = \exp(-\beta \|\mathbf{x} - \mathbf{t}\|^2)$ for $\beta > 0$ the data $\mathbf{x}, \mathbf{t} \in \mathbb{R}^n$ (It has infinite dimensional feature map)
- ANOVA kernel: $k(\mathbf{x}, \mathbf{t}) = \prod_{i=1}^n (1 + x_i t_i)$

Remark 13. Consider the following polynomial kernel as we have:

$$\sum_{i=1}^d \frac{a^i}{i!} (\mathbf{x}^T \mathbf{t})^i$$

Suppose we have $r = \infty$, this can converge uniformly to $\exp(a\mathbf{x}^T \mathbf{t})$ showing that it is a kernel, where if $n = 1$, the feature map is:

$$\phi = \left(1, \sqrt{2}x, \sqrt{\frac{a}{2}}x^2, \sqrt{\frac{a^3}{6}}x^3, \dots\right) = \left(\sqrt{\frac{a^i}{i!}} : i \in \mathbb{N}\right)$$

Definition 2.10. (Transition Invariance/Radial Kernel) We say that a kernel $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is:

- *Transition Invariance:* If the kernel has the form:

$$k(\mathbf{x}, \mathbf{t}) = H(\mathbf{x} - \mathbf{t})$$

for all $\mathbf{x}, \mathbf{t} \in \mathbb{R}^d$ where $H : \mathbb{R}^d \rightarrow \mathbb{R}$ is a differentiable function.

- *Radial,* if kernel has the form:

$$k(\mathbf{x}, \mathbf{t}) = h(\|\mathbf{x} - \mathbf{t}\|)$$

for all $\mathbf{x}, \mathbf{t} \in \mathbb{R}^d$ where $h : [0, \infty) \rightarrow [0, \theta]$ is the differentiable function.

Remark 14. The important example of a radial kernel in the Gaussian kernel as we have:

$$k(\mathbf{x}, \mathbf{t}) = \exp(-\beta \|\mathbf{x} - \mathbf{t}\|^2)$$

which is a product of 2 kernel as $k(\mathbf{x}, \mathbf{t}) = \exp(-\beta(\mathbf{x}^T \mathbf{x} + \mathbf{t}^T \mathbf{t})) \exp(2\beta \mathbf{x}^T \mathbf{t})$

Remark 15. (Ridge Regression with Feature Map) Given the dataset $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $\mathbf{y} \in \mathbb{R}^{m \times 1}$. Starting with the basis function ϕ_1, \dots, ϕ_N where $\phi_i : \mathbb{R}^n \rightarrow \mathbb{R}$ with the map:

$$\Phi = \begin{bmatrix} \phi_1(\mathbf{x}_1) & \dots & \phi_N(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_m) & \dots & \phi_N(\mathbf{x}_m) \end{bmatrix} \in \mathbb{R}^{m \times N}$$

We have the regression coefficient as we have $\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I}_N)^{-1} \Phi^T \mathbf{y}$

Remark 16. (Kernel Ridge Regression) Given the same setting, a kernel function $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$, where the kernel matrix is given by:

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \in \mathbb{R}^{m \times m}$$

Regression coefficient is then given by $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_m)^{-1} \mathbf{y}$ as the function is:

$$\hat{y}(\mathbf{x}) = \sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

3 Support Vector Machine

3.1 Forming Problems

Definition 3.1. (Seperating Hyperplane) Let the dataset be $S = \{(x_i, y_i)\}_{i=1}^m \in \mathbb{R}^n \times \{-1, 1\}$. The hyperplane is the set such that:

$$\mathcal{H}_{w,b} = \{x \in \mathbb{R}^n : w^T x + b = 0\}$$

Definition 3.2. (Linearly Separatable) The data are linearly separatable if there exists $w \in \mathbb{R}^n$ and $b \in \mathbb{R}$ such that:

$$y_i(w^T x_i + b) > 0$$

for $i = 1, \dots, m$, which we call $\mathcal{H}_{w,b}$ a separating hyperplane. Note that it is a strict inequality.

Proposition 3.1. (Finding A distance from Plane) If $\mathcal{H}_{w,b}$ is a hyperplane, we also define the distance from a point x to be:

$$\frac{w^T x + b}{\|w\|}$$

Proof. We consider the projection from the point x to $\mathcal{H}_{w,b}$ as we have:

$$p = x - \frac{w(b + w^T x)}{\|w\|^2}$$

To show that p is indeed a projection:

- We will have to show that p is on hyperplane

$$w^T p + b = w^T x - \frac{w^T w(b + w^T x)}{\|w\|^2} + b = 0$$

- $x - p$ is orthogonal to $p - x'$ where x' is any point from on the hyperplane:

$$\begin{aligned} (p - x)^T (p - x') &= \left\langle -\frac{w(b + w^T x)}{\|w\|^2}, p - x' \right\rangle \\ &= \left\langle -\frac{w(b + w^T x)}{\|w\|^2}, x - \frac{w(b + w^T x)}{\|w\|^2} - x' \right\rangle \\ &= \left\langle -\frac{w(b + w^T x)}{\|w\|^2}, x - x' \right\rangle + \left\langle \frac{w(b + w^T x)}{\|w\|^2}, \frac{w(b + w^T x)}{\|w\|^2} \right\rangle \\ &= \left\langle -\frac{w(b + w^T x)}{\|w\|^2}, x - x' \right\rangle + \frac{\|w\|^2 (b + w^T x)^2}{\|w\|^4} \\ &= -\frac{b + w^T x}{\|w\|^2} \langle w, x - x' \rangle + \frac{(b + w^T x)^2}{\|w\|^2} \\ &= -\frac{(b + w^T x)(w^T x - w^T x')}{\|w\|^2} \langle w, x - x' \rangle + \frac{(b + w^T x)^2}{\|w\|^2} \\ &= -\frac{b(w^T x) - b(w^T x') + (w^T x)^2 - (w^T x)(w^T x')}{\|w\|^2} + \frac{(b + w^T x)^2}{\|w\|^2} \\ &= -\frac{b(w^T x) + b^2 + (w^T x)^2 + (w^T x)b}{\|w\|^2} + \frac{(b + w^T x)^2}{\|w\|^2} = 0 \end{aligned}$$

Please note that $w^T x' + b = 0$.

Now, we are left to find the distance between \mathbf{p} and \mathbf{x} , which we can find it to be:

$$\sqrt{(\mathbf{p} - \mathbf{x})^T(\mathbf{p} - \mathbf{x})} = \sqrt{\left\langle \frac{\mathbf{w}(b + \mathbf{w}^T \mathbf{x})}{\|\mathbf{w}\|^2}, \frac{\mathbf{w}(b + \mathbf{w}^T \mathbf{x})}{\|\mathbf{w}\|^2} \right\rangle} = \frac{|b + \mathbf{x}^T \mathbf{w}|}{\|\mathbf{w}\|}$$

Thus complete the proof. \square

Definition 3.3. (Margin) As we have the distance from a point \mathbf{x} to the plane $\mathcal{H}_{\mathbf{w},b}$ to be $\rho_{\mathbf{x}}(\mathbf{w}, b)$. If $\mathcal{H}_{\mathbf{w},b}$ separates the training set S , we define a margin as:

$$\rho_S(\mathbf{w}, b) = \min_{i \in [m]} \rho_{\mathbf{x}_i}(\mathbf{w}, b)$$

Definition 3.4. (Optimal Separating Hyper-Planes) We want to find the weight and bias of a separating hyperplane such that the the margin is maximized :

$$\rho(S) = \max_{\mathbf{w}, b} \min_{i \in [m]} \left\{ \frac{y_i(\mathbf{w}^T \mathbf{x}_i + b)}{\|\mathbf{w}\|} : y_j(\mathbf{w}^T \mathbf{x}_j + b) > 0 \text{ for } j \in [m] \right\}$$

Furthermore, to get the unique \mathbf{w}, b , we may consider 2 choices:

- Set $\|\mathbf{w}\| = 1$, so $\rho_{\mathbf{x}}(\mathbf{w}, b) = |\mathbf{w}^T \mathbf{x} + b|$ and so:

$$\rho_S = \min_{i \in [m]} y_i(\mathbf{w}^T \mathbf{x}_i + b)$$

- Choose $\|\mathbf{w}\|$ such that $\rho_S(\mathbf{w}, b) = 1/\|\mathbf{w}\|$ or:

$$\min_{i \in [m]} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$$

We will consider the second case.

Proposition 3.2. *The optimal separating hyperplane is equivalent to following optimization problem:*

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{aligned}$$

for $\mathbf{w} \in \mathbb{R}^n$. The quantity $1/\|\mathbf{w}\|$ is the margin of optimal separating hyperplane.

Proof. We have following the second case:

$$\begin{aligned} \rho(S) &= \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} : \min_{j \in [m]} \{y_j(\mathbf{w}^T \mathbf{x}_j + b)\} = 0, y_k(\mathbf{w}^T \mathbf{x}_k + b) > 0 \text{ for } k \in [m] \right\} \\ &= \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} : \{y_j(\mathbf{w}^T \mathbf{x}_j + b)\} \geq 1 \right\} = \frac{1}{\min_{\mathbf{w}, b} \{\|\mathbf{w}\| : \{y_j(\mathbf{w}^T \mathbf{x}_j + b)\} \geq 1\}} \end{aligned}$$

\square

Proposition 3.3. *To minimize a differentiable convex function $f(\mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}$ subjected to linear inequality $\mathbf{A}\mathbf{z} \leq \mathbf{c}$. We may solve the problem with Lagrangian:*

$$L(\mathbf{x}, \boldsymbol{\alpha}) = f(\mathbf{x}) - \boldsymbol{\alpha}^T(\mathbf{A}\mathbf{x} - \mathbf{c})$$

If the optimization problem is feasible that is $\{\mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{c}\} \neq \emptyset$, we can show that:

$$\max_{\boldsymbol{\alpha} \geq \mathbf{0}} \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}) = \min_{\mathbf{x}} f(\mathbf{x}) \text{ s.t } \mathbf{A}\mathbf{x} \leq \mathbf{c}$$

And there is a necessary and sufficient condition called KKT for a solution $(\boldsymbol{\alpha}^, \mathbf{z}^*)$:*

- $\mathbf{A}\mathbf{x}^* \leq \mathbf{c}$
- $\boldsymbol{\alpha}^* \geq \mathbf{0}$
- $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha})|_{\mathbf{x}^*} = \mathbf{0}$
- $(\mathbf{A}\mathbf{x}^* - \mathbf{c})_i \alpha_i^* = 0_i$ for $i \in [m]$

Proposition 3.4. *The dual form for the SVM is:*

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & -\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha} + \sum_{i=1}^m \alpha_i \\ \text{s.t} \quad & \sum_{i=1}^m y_i \alpha_i = 0 \text{ for } i \in [m] \\ & \alpha_i \geq 0 \end{aligned}$$

where $\mathbf{A} = (y_i y_j \mathbf{x}_i^T \mathbf{x}_j : i, j \in [m])$. The solution to the primal problem is:

$$\mathbf{w}^* = \sum_{i=1}^m \alpha_i^* y_i \mathbf{x}_i$$

as the weight is the linear combination of the data. Finally the variable b^* can be determine by find the weight \mathbf{x}_j that satisfies the condition:

$$y_i((\mathbf{w}^*)^T \mathbf{x}_i + b) - 1 = 0$$

Then we bias can be found by rearrange as we have $b^* = y_i - (\mathbf{w}^*)^T \mathbf{x}_j$. The point that satisfies this conditon is called support vector.

Proof. We consider the Lagragian to be:

$$L(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1]$$

where $\alpha_i \geq 0$ is Lagragian multipler. Let's minimize L over \mathbf{w} and b and maximized over $\boldsymbol{\alpha}$ with $\boldsymbol{\alpha} \geq \mathbf{0}$. We can see that the partial derivative is:

$$\begin{aligned} \frac{\partial L}{\partial b} &= -\sum_{i=1}^m y_i \alpha_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \end{aligned}$$

Now, we can see that the optimal weight will have the linear combination term. Let's plugging this back into Lagragian and we have:

$$\frac{1}{2} \underbrace{\mathbf{w}^T \mathbf{w}}_{\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha}} - \underbrace{\sum_{i=1}^m \alpha_i y_i \mathbf{w}^T \mathbf{x}_i}_{\boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha}} - b \underbrace{\sum_{i=1}^m \alpha_i y_i}_0 + \sum_{i=1}^m \alpha_i$$

□

Remark 17. The new point \mathbf{x} can be classified as:

$$\text{sign} \left(\sum_{i=1}^m y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x}_i + b^* \right)$$

One can show that the expected generalization error of SVM trained on $m - 1$ sample is bounded by n_{sv}/m , where n_{sv} is the number of support vector.

Remark 18. (Linear Non-Separatable Case) We would like to minimize the following objective function:

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m V_{\text{mc}}(y_i, \mathbf{w}^T \mathbf{x}_i + b)$$

as we have $V_{\text{mc}}(y, \hat{y}) = \mathbb{I}[y = \text{sign}(\hat{y})]$ but it is NP-Hard and so we will have to convexify the problem by consider the hinge loss, instead:

$$V_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - h\hat{y})$$

This will gives us the convex optimization.

Proposition 3.5. *The hinge loss can be reformulated using the slack variable and gives us the following optimization problem:*

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \text{ for } i = 1, \dots, m \end{aligned}$$

This would in turn, gives us the following dual problem:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & -\frac{1}{2} \boldsymbol{\alpha}^T \mathbf{A} \boldsymbol{\alpha} + \sum_{i=1}^m \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^m y_i \alpha_i = 0 \text{ for } i \in [m] \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

We will consider the implication of KKT conditon afterward.

Proof. We now have the following Lagragian to be:

$$L(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1] - \sum_{i=1}^m \beta_i \xi_i$$

where $\alpha_i, \beta_i \geq 0$ are Lagragian multipler. We minimize L over $(\mathbf{w}, \boldsymbol{\xi}, b)$ and maximize L with respected to the variables as:

$$\begin{aligned} \frac{\partial L}{\partial b} &= -\sum_{i=1}^m y_i \alpha_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} &= \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \\ \frac{\partial L}{\partial \xi_i} &= c - \alpha_i - \beta_i = 0 \implies 0 \leq \alpha_i \leq C \end{aligned}$$

Plugging this back gives us the dual form. Please note that both $\alpha_i, \beta_i \geq 0$ □

Remark 19. (Interpretation of The Results) The dual problem is similar to the eariler linear separatable case, as we have additional box constraint. The weight is given as:

$$\mathbf{w}^* = \sum_{i=1}^m \alpha_i^* y_i \mathbf{x}_i$$

where \mathbf{b}^* is the same. For a new KKT conditon, we have:

$$\begin{aligned} \alpha_i^* (y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* - 1 + \xi_i^*) &= 0 \\ (C - \alpha_i^*) \xi_i^* &= 0 \end{aligned}$$

where the second equation follows from $\beta_i^* = C - \alpha_i^*$. There are difference points to consider:

- $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* > 1$ implies that $\alpha_i^* = 0$ where the point isn't support vector.
- $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* < 1$ implies that $\alpha_i^* = C$ where the point is a support vector slack ξ_i^* outlier.
- $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* = 1$ implies that $\alpha_i^* \in [0, C]$ and if $\alpha_i^* > 0$, it is a support vector on a margin.

On the otherhand, we have:

- $\alpha_i^* = 0$ then we have $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* \geq 1$ and $\xi_i^* = 0$
- $\alpha_i^* \in (0, C)$ then we have $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* = 1$ and $\xi_i^* = 0$
- $\alpha_i^* = C$ then we have $y_i(\mathbf{w}^*)^T \mathbf{x}_i + b^* \leq 1$ and $\xi_i^* \geq 0$

Remark 20. The role of parameter C is that:

- The parameter C controls the trade-off between $\|\mathbf{w}\|^2$ and the training error $\sum_{i=1}^m \xi_i$
- The value of α_i^* is piecewise quadratic of C
- C is selected by minimizing leave-one-out (LOO) cross-validation error.

To compute the LOO error, we need to retrain the SVM no more than the number of support vector making it fast to train. One can observe that we can use the n_{sv}/m as an upper bound on LOO error.

Definition 3.5. (Kernelized SVM) Given the feature map $\phi(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{W}$, we can replace \mathbf{x} with $\phi(\mathbf{x})$ and $\mathbf{x}^T \mathbf{t}$ by $\langle \phi(\mathbf{x}), \phi(\mathbf{t}) \rangle$. The result function is:

$$f(\mathbf{x}) = \sum_{i=1}^m y_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b$$

The parameter can be found using the matrix $\mathbf{A} = (y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) : i, j \in [m])$ and the new point is classified the same.

Remark 21. (Connection to the Regularization) SVM formulation is equivalent to the following problem:

$$\mathcal{E}_\lambda = \sum_{i=1}^m \max \left(1 - y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b), 0 \right) + \lambda \|\mathbf{w}\|^2$$

where we set $\lambda = 1/(2C)$ and so we have:

$$\begin{aligned} & \min_{\mathbf{w}, b, \xi} \left\{ C \sum_{i=1}^m \xi_i + \frac{1}{2} \|\mathbf{w}\|^2 : y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0 \right\} \\ &= \min_{\mathbf{w}, b} \left\{ \min_{\xi} \left\{ C \sum_{i=1}^m \xi_i + \frac{1}{2} \|\mathbf{w}\|^2 : y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0 \right\} \right\} \\ &= \min_{\mathbf{w}, b} \left\{ C \sum_{i=1}^m \left(1 - y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b), 0 \right) + \frac{1}{2} \|\mathbf{w}\|^2 \right\} = C \mathcal{E}_{1/(2C)}(\mathbf{w}, b) \end{aligned}$$

Remark 22. (SVM for Regression) If we have the regression for the SVM, then we use the following loss:

$$|y - f(\mathbf{x})|_\varepsilon = \max(|y - f(\mathbf{x})| - \varepsilon, 0)$$

This would gives the following optimization problem:

$$\begin{aligned} \min \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{x}_i + b - y_i \leq \varepsilon + \xi_i \\ & y_i - \mathbf{w}^T \mathbf{x}_i - b \leq \varepsilon + \xi_i^* \\ & \xi_i, \xi_i^* \geq 0 \text{ for } i \in [m] \end{aligned}$$

Please note that the loss function is scale sensitive as the error below certain. This gives the sparse solution. One can use decompositive to solve all of the KKT problems.

4 Tree Based and Ensemble Model

4.1 Tree Based Method

Definition 4.1. (Tree Method) We are interesting to partition the input space into retangles and fit simple model in each one; for example, we have the function:

$$f(\mathbf{x}) = \sum_{p=1}^P c_p \mathbb{I}[\mathbf{x} \in R_p]$$

Where we hve the following:

- We partition the input space with hyper-retangle R_1, R_2, \dots, R_P where: $\bigcup_{p=1}^P R_p = \mathcal{X}$ and $R_a \cap R_b = \emptyset$ if $a \neq b$
- $\{c_p\}_{p=1}^P$ is some real parameter with a natural choice to be:

$$c_p = \text{avg}(y_i | \mathbf{x}_i \in R_p) = \frac{\sum_{i=1}^m y_i \mathbb{I}[\mathbf{x}_i \in R_p]}{\sum_{i=1}^m \mathbb{I}[\mathbf{x}_i \in R_p]}$$

We are interested to solving the following optimization problem:

$$\min_{R_1, \dots, R_P} \left\{ \sum_{i=1}^m \left(y_i - \sum_{p=1}^P \text{avg}(y_i | \mathbf{x}_i \in R_p) \mathbb{I}[\mathbf{x}_i \in R_p] \right)^2 \right\}$$

Definition 4.2. (Heuristic Search) It seem to be intractable, so we need heuristic approach. Let's find the way to split the tree. Define a pair of axis parallel half-spaces:

$$R_1(j, s) = \{\mathbf{x} | x_j \leq s\} \quad R_2(j, s) = \{\mathbf{x} | x_j > s\}$$

Then we search for optimal values j^* and s^* , which solves the problem:

$$\min_{j, s} \left\{ \min_{c_1} \sum_{\mathbf{x}_i \in R_1(j, s)} (y_i - c_1(\mathbf{x}_i))^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j, s)} (y_i - c_2(\mathbf{x}_i))^2 \right\}$$

The inner minimizer is solved by:

$$c_1^* = \text{avg}(y_i | \mathbf{x}_i \in R_1(j, s)) \quad c_2^* = \text{avg}(y_i | \mathbf{x}_i \in R_2(j, s))$$

For each splitting variable j , the search for best split at point s can be don by $\mathcal{O}(m)$ computation. Thus, the problem is solved in $\mathcal{O}(nm)$ computation. The decision tree can be solved by repeatedly splitting the tree branches.

Remark 23. (Overfitting) If we keep repeating the heuristic search process, we will overfit the data. There are several ways to fix this:

- The following the split only if it decreases the empirical error more than the threshold. However, this might be the best as we might find split below a bad mode.
- We might consider the maximal depth of split tree is reached. This could leads to an underfitting or overfitting. We need to look at the data to determine the size of tree.

Remark 24. (Solving Overfitting) We choose the tree adapting from the data. We grows the large tree \hat{T} (stopping when the maximum number of data is assigned at each node). Now consider the prune the tree with cost complexity pruning i.e looks for subtree $T_\lambda \subseteq \hat{T}$ that minimizes:

$$C_\lambda(T) = \sum_{p=1}^{|T|} m_p Q_p(T) + \lambda |T|$$

where T is the subtree of \hat{T} , where we have:

- p runs over leaf nodes of T (a subset of the nodes of \hat{T})
- m_p is the number of data point assigned to node p
- Q_p is the training error given as:

$$Q_p = \frac{1}{m_p} \sum_{\mathbf{x}_i \in R_p} (y_i - c_p)^2$$

At the first term in C_λ is the training error.

One can show that there is a unique $T_\lambda \subseteq \hat{T}$, with minimize C_λ , while a good value of λ can be found by cross-validation.

Definition 4.3. (Weakest Link Pruning) We successively collapse the internal nodes that produces the smallest per node increase in:

$$\sum_{p=1}^{|T|} m_p Q_p(T)$$

We continue until the root the tree is produce. As now, we have a list of pruned trees. We can search along this list for the one that minimizes the objective C_λ , and one can show that T_λ is in the produced list of subtree, hence the algorithm gives the optimal solution.

Definition 4.4. (Classification Tree) When the output is a categorical variable, we use the same algorithm above with 2 important modification:

- For each region R_n , we define the empirical class probability, as we have:

$$p_{nk} = \frac{1}{m_p} \sum_{(\mathbf{x}_i, y_i) \in R_n} \mathbb{I}[y_i = k]$$

- We classify an input which falls in region n in the class with new probability as we have:

$$f(\mathbf{x}) = \arg \max_{k \in \{1, \dots, K\}} \sum_{n=1}^N p_{nk} \mathbb{I}[\mathbf{x} \in R_n]$$

Definition 4.5. (Impurity) We consider the training error $Q_p(T)$ to be called impurity, which in can be one of these values:

- *Misclassification Error*: $1 - p_{pk(n)}$ where $k(n) = \arg \max_{k \in \{1, \dots, K\}} p_{nk}$
- *Gini-Index*: $\sum_k p_{pk}(1 - p_{pk})$
- *Cross-Entropy*: $\sum_k p_{pk} \log(1/p_{pk})$

The cross-entropy or gini-index are used to growing the tree, while the misclassification error are often used to prune the tree.

4.2 Ensemble Methods + Bagging

Theorem 4.1. (Chernoff-Bound) Let X_1, X_2, \dots, X_n be independent random variable. Assuming $0 \leq x_i \leq 1$. We denote the $X = \sum_{i=1}^n X_i$ and $\mu = \mathbb{E}[X] = \sum_{i=1}^n \mathbb{E}[X_i]$, then for all $0 \leq k \leq \mu$:

$$\mathbb{P}(X \leq k) \leq \exp\left(-\frac{(\mu - k)^2}{2\mu}\right)$$

Remark 25. (Motivation - Wisdom of the Crowd) A single individual might often wrong but the crowd majority may often be corrected. Suppose each individual in the crowd $h_1, h_2, \dots, h_{2T+1}$ of the size $2T + 1$ predicts the outcome correctly with probability $1/2 + \gamma$ independent from each other. We consider the vote of the crowd to be:

$$H_T = \text{sgn}\left(\sum_{t=1}^{2T+1} h_t\right)$$

The probability of H_T being wrong is given as:

$$\mathbb{P}(H_T \text{ is wrong}) = \sum_{i=1}^T \binom{2T+1}{i} \left(\frac{1}{2} + \gamma\right)^i \left(\frac{1}{2} - \gamma\right)^{2T+1-i}$$

We simplify the above using a Chernoff bound. We let $X_1, \dots, X_i, \dots, X_n$ be Bernoulli random variable where $X_i = 1$ if voter i is correct and 0 otherwise. Taking $k = T$ and $n = 2T + 1$ thus:

$$\mu = (2T + 1) \left(\frac{1}{2} + \gamma\right) = T + \frac{1}{2} + 2T\gamma + \gamma$$

Now, we substitute the bound:

$$\begin{aligned} \mathbb{P}(H_T \text{ is wrong}) &\leq \exp\left(-\frac{(\mu - T)^2}{2\mu}\right) \\ &= \exp\left(-\frac{(1/2 + 2T\gamma + \gamma)^2}{2(T + 1/2 + 2T\gamma + \gamma)}\right) \\ &\leq \exp\left(-\frac{4T^2\gamma^2}{5T}\right) = \exp\left(-\frac{4\gamma^2}{5}T\right) \end{aligned}$$

The bound may be too crude but the probability of getting wrong, exponentially decays to zero.

Definition 4.6. (Bagging Algorithm) The idea of bagging algorithm is to reduce the variance of a classifier by having many variances of the classifier and then voting. We have the following algorithm:

- Training data: $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\} \subset \mathbb{R}^d \times \{-1, 1\}$
- Ensemble of size T
- Resample dataset of size M
- Classifier function $h_S(\mathbf{x})$

This leads to the following pseudocode:

Algorithm 1 Bagging Algorithm

- 1: **for** $t = 1, 2, \dots, T$ **do**
- 2: $S[t] = M$ examples sampled with replacement from S
- 3: **end for**
- 4: **Return:** We perform the following prediction:

$$H(\mathbf{x}) = \text{sgm}\left(\sum_{t=1}^T h_{S[t]}(\mathbf{x})\right)$$

We may set M to be m .

Remark 26. If we set $M = m$, we can find the number of unique example from S are in bag $S(t)$. The probability that a particular example doesn't appear in the bag is $(1 - 1/m)^m$, and please note that:

$$\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^m = \frac{1}{e} \approx 0.368..$$

so there will be around 63% examples in each dataset $S[t]$.

Definition 4.7. (Random Forest) We observe the wisdom of the crowds argument. We can build a tree using a subset of size k features, which is usually \sqrt{d} or $\log d$.

4.3 Boosting

Remark 27. (Concept of Boosting) Some of the problem is easy to find the “rule of thumb” that is usually correct. It is hard to find accurate prediction rule. To boosting algorithm is given by:

- Create a computer program for deriving rough rule of thumb.
- We can shoow a rule of thumb to fit a subset of example.
- Repeat T times.
- Combined the classifier by weighted majority votes.

There are two concerns: How do we choose the subset of examples ? At each round as we want to concentrate on the hardest example. How do we combine the weak learner ? This can be done by weighted majority.

Definition 4.8. (Notation Used in Boosting) We have the following variables, as we have:

- $D_t(i)$: Weight on example i at time t when $\sum_{i=1}^m D_t(i) = 1$
- α_t : Weight on weak learner t where $\alpha_t \in \mathbb{R}$
- $h_t(\cdot) : \mathbb{R}^d \rightarrow \{-1, +1\}$: Weak learner that is generated at time t .
- $f(\cdot)$: Weighted on weak learner. $\sum_{t=1}^T \alpha_t h_t(\mathbf{x})$
- $H(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$: Final classifier.
- ε_t : Weight error of weak learner $h_t(\cdot)$ at time t :

$$\varepsilon_t = \sum_{i=1}^m D_t(i) \mathbb{I}[h_t(\mathbf{x}_i) \neq y_i]$$

- Weak learning will generate the output:

$$D_t(1), \dots, D_t(m), (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$$

The weak-learner will output a weaker learner $h_t(\cdot)$ such that $\varepsilon_t < 1/2$

Definition 4.9. (Adaboost Algorithm) We have the following pseudocode for the adaboost this is shown in the pseudocode [2](#).

Algorithm 2 Adaboost

- 1: **Input:** Training set $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$
- 2: **Initialize:** $D_1(1) = \dots = D_1(m) = 1/m$
- 3: **for** $i = 1, 2, \dots, T$ **do**
- 4: Fit the classifier $h_t : \mathbb{R}^d \rightarrow \{-1, 1\}$ using a distribution D_t
- 5: Choose $\alpha_t \in \mathbb{R}$:

$$\alpha_t = \frac{1}{2} \log \frac{1 - \varepsilon_t}{\varepsilon_t}$$

- 6: Update for each $i \in [m]$, where Z_t is normalization factor:

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))}{Z_t}$$

- 7: **end for**

- 8: **Return:** Classifier is given as:

$$H(\mathbf{x}) = \text{sgn} \left(\sum_{t=1}^T \alpha_t h_t(\mathbf{x}) \right)$$

Typically $\varepsilon_t \leq 0.5$ hence $\alpha_t \geq 0$. Thus f is a linear combination of h_t with weights controlled by training error. The basic intuition for the adaboost assign a larger weight are assigned to hard examples, hence the weak learner will focus on those example.

Theorem 4.2. *Given a training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ and assume that each iteration of Adaboost the weak learner returns a hypothesis with a weighted error $1/2 - \gamma \geq \varepsilon_t$, then training error of the output hypothesis is at most:*

$$\frac{1}{m} \sum_{i=1}^m \mathbb{I}[H(\mathbf{x}_i) \neq y_i] \leq \exp(-2\gamma^2 T)$$

Proof. Please note that the training error is bounded as:

$$\frac{1}{m} \sum_{i=1}^m \mathbb{I}[H(\mathbf{x}_i) \neq y_i] \leq \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(\mathbf{x}_i))$$

where $f = \sum_t \alpha_t h_t$ so that $H(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$. The inequality follows from $H(\mathbf{x}_i) \neq y_i$ implies that $\exp(-y_i f(\mathbf{x}_i)) \geq 1$. Now consider the definition of D_t where, recursively:

$$D_{T+1}(i) = \frac{1}{m} \frac{\prod_{t=1}^T \exp(-\alpha_t y_i h_t(\mathbf{x}_i))}{\prod_{t=1}^T Z_t}$$

We can expand this equation, where we have:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \exp(-y_i f(\mathbf{x}_i)) &= \frac{1}{m} \sum_{i=1}^m \exp \left(-y_i \sum_{t=1}^T \alpha_t h_t(\mathbf{x}_i) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \prod_{t=1}^T \exp(-y_i \alpha_t h_t(\mathbf{x}_i)) \\ &= \sum_{i=1}^m D_{T+1}(i) \prod_{t=1}^T Z_t = \prod_{t=1}^T Z_t \end{aligned}$$

If at each iteration, we choose α_t and h_t by minimizing Z_t , the final training error of H will be reduced most rapidly. Recall that:

$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha_t y_i h_t(\mathbf{x}_i))$$

Using the fact that Z_t is a binary, we have that:

$$\begin{aligned} Z_t &= \exp(\alpha_t) \sum_{i: y_i \neq h_t(\mathbf{x}_i)} D_t(i) + \exp(-\alpha_t) \sum_{i: y_i = h_t(\mathbf{x}_i)} D_t(i) \\ &= \varepsilon_t \exp(\alpha_t) + (1 - \varepsilon_t) \exp(-\alpha_t) \end{aligned}$$

Setting the derivative of Z_t to zero with respect to α_t , which gives us the weight:

$$\alpha_t = \frac{1}{2} \log \frac{1 - \varepsilon_t}{\varepsilon_t}$$

Placing α_t to the value Z_t , and we have:

$$\begin{aligned} Z_t &= \varepsilon_t \exp(\alpha_t) + (1 - \varepsilon_t) \exp(-\alpha_t) \\ &= 2\sqrt{\varepsilon_t(1 - \varepsilon_t)} = \sqrt{1 - 4\gamma_t^2} \end{aligned}$$

Please note that $\gamma_t = 1/2 - \varepsilon_t$. Hence we have:

$$\frac{1}{m} \sum_{i=1}^m \mathbb{I}[H(\mathbf{x}_i) \neq y_i] \leq \prod_{t=1}^T Z_t = \prod_{t=1}^T \sqrt{1 - 4\gamma_t^2} \leq \exp\left(-2 \sum_{t=1}^T \gamma_t^2\right)$$

The final inequality use the fact that $1 - x \leq \exp(x)$. If each weak classifier is slightly better than random guessing, the training drops exponentially fast. \square

Remark 28. (Derivation of Adaboost) The boosting can be seen as a greedy way to solve problem:

$$\min \left\{ \sum_{i=1}^m V\left(y_i, \sum_{i=1}^T \alpha_t h_t(\mathbf{x}_i)\right) : \alpha_1, \dots, \alpha_T \in \mathbb{R}^T, h_1, \dots, h_T \in \mathcal{H}^T \right\}$$

where \mathcal{H} is hypothesis class which contains the weaker learner and the loss function is exponential for instance $V(y, \hat{y}) = \exp(-y\hat{y})$. At each iteration, a new basis function is added to the current basis expansion $f^{(t-1)} = \sum_{s=1}^{t-1} \alpha_s h_s$, which we have:

$$(\alpha_t, h_t) = \arg \min_{\alpha_t, h_t} \sum_{i=1}^m V(y_i, f^{(t-1)}(\mathbf{x}_i) + \alpha_t h_t(\mathbf{x}_i))$$

unlike the decision tree, where each iteration in previous basis is re-adjusted. In statistics literature, this kind of model is called stagewise additive model. To derive the adaboost, substitute $V(y, \hat{y}) = \exp(-y\hat{y})$ and we consider the following optimization problem:

$$\min_{\alpha_t, h_t} \sum_{i=1}^m \exp\left(-y_i \left(f^{(t-1)}(\mathbf{x}_i) + \alpha_t h_t(\mathbf{x}_i)\right)\right)$$

We define $\mathcal{D}_t(i) = \exp(-y_i f^{(t-1)}(\mathbf{x}_i))$ as we have:

$$\min_{\alpha_t, h_t} \sum_{i=1}^m \mathcal{D}_t(i) \exp(-\alpha_t h_t(\mathbf{x}_i) y_i)$$

We can see that the This equation can be rewritten as:

$$\begin{aligned} &\min_{\alpha_t, h_t} \left(\exp(\alpha_t) \sum_{i: y_i \neq h_t(\mathbf{x}_i)} \mathcal{D}_t(i) + \exp(-\alpha_t) \sum_{i: y_i = h_t(\mathbf{x}_i)} \mathcal{D}_t(i) \right) \\ &= \min_{\alpha_t, h_t} \left((e^{\alpha_t} - e^{-\alpha_t}) \sum_{i=1}^m \mathcal{D}_t(i) \mathbb{I}[y_i \neq h_t(\mathbf{x}_i)] + e^{-\alpha_t} \sum_{i=1}^m \mathcal{D}_t(i) \right) \end{aligned}$$

This is similar to the adaboost, which we have: h_t minimizes the weight misclassification error weight by \mathcal{D}_t that is is propotional to adaboost D_t . Finally, minimization of α_t is the same as adaboost.

Remark 29. (Classification and Regression) In the typical setup of classification as we have:

$$\min_{f \in \mathcal{F}} \sum_{i=1}^m V(y_i, f(\mathbf{x}_i)) + \lambda \text{ complexity}(f)$$

There are some problems with classification as we use the exponential loss. To make the class of function \mathcal{F} both rich and smooth, we have the function f that maps to \mathbb{R} rather than $\{-1, 1\}$ then predict the sign. We have the typical loss function, where we have for $y \in \{-1, +1\}$:

- Misclassification Loss: $V_{\text{mc}}(y, \hat{y}) = \mathbb{I}[y \neq \text{sgn}(\hat{y})]$. It isn't continuous.
- Hinge Loss: $V_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - y\hat{y})$. It punishes the negative margin but not positive margin, but it isn't differentiable everywhere.
- Square Loss: $V_{\text{sq}}(y, \hat{y}) = (y - \hat{y})^2$. It unnecessarily punishes predicting with increasing positive margin.
- Exponential Loss: $V_{\text{exp}}(y, \hat{y}) = \exp(-y\hat{y})$. It punishes negative margins and promotes large positive margins.

Thus the exponential loss is chosen.

5 Online Learning

5.1 Introduction

Definition 5.1. (Online Learning with Expert Advice) There exists an online sequence of data:

$$S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \{0, 1\}^n \times \{0, 1\}$$

The vector \mathbf{x}_t is the set of predictions of n experts at time t , which we aim to predict y_t . We would like to find an algorithm that tried to combine the predictions \mathbf{x}_t of the n experts to predict \hat{y}_t , an estimate of y_t . The loss of master algorithm A on sequence S :

$$L_A(S) = \sum_{t=1}^m |y_t - \hat{y}_t|$$

We want to find an algorithm with a small loss.

Definition 5.2. (Regret) Recall the loss function $L_A(S)$ and we let:

$$L_i(S) = \sum_{t=1}^m |y_t - x_{t,i}|$$

being the loss of i -th expert E_i . The aim of our algorithm should be that the bound of the form, such that for all sequence S :

$$L_A(S) \leq a \min_i L_i(S) + b \log(n)$$

where a, b are small constants. This is known as regret as it is the loss of objective related to the best expert.

Definition 5.3. (Halving Algorithm) Suppose that there is an expert that is consistent (gives correct answer), we can perform the search on this consistent expert, in which we will have the correct prediction:

- If a mistake is made, the number of consistent experts is (at least) halved.
- For any sequence with a consistent expert, the Halving algorithm made less than or equal to $\log_2 n$ mistakes.

5.2 Learning from Expert Algorithm

Definition 5.4. (Weighted Majority) This algorithm for non-consistent expert, which we can perform the prediction with larger scale. We have weight of the wrong expert is multiplied by $\beta \in [0, 1)$.

Theorem 5.1. *The number of mistake of master algorithm M , with $\beta = 1/2$ is given by:*

$$M \leq 2.63 \min_i M_i + 2.63 \ln n$$

where M_i is the number of mistakes of expert E_i

Proof. We have the following quantities:

- $M_{t,i}$ is the number of mistake of the expert i , E_i at the start of trial t .
- $w_{t,i} = \beta^{M_{t,i}}$ weights of E_i at the begin of trial t .
- Please note that $w_{1,i} = 1$, and $W_t = \sum_{i=1}^n w_{t,i}$ is the total weight at trial t .

It is clear that the total weight of the minority is when it is less than $1/2 W_t$, but the total weight of the majority is when it is more than $1/2 W_t$. There are 2 scenarios, which we have:

- If no mistake, the minority expert weight is multiplied by β as we have (Trivial Bound): $W_{t+1} \leq 1 \cdot W_t$
- If there is a mistake, the majority expert weights are multiplied by β as:

$$\begin{aligned} W_{t+1} &= \text{Minority} + \beta \text{Majority} \\ &\leq \frac{1}{2} W_t + \beta \text{Majority} \\ &\leq \frac{1}{2} W_t + \beta \frac{1}{2} W_t \leq \frac{1+\beta}{2} W_t \end{aligned}$$

The third inequality comes from the fact that the majority is at least $1/2 W_t$, making the upperbound tighter.

This gives us:

$$\left(\frac{1+\beta}{2}\right)^M W_1 \geq W_{m+1} = \sum_{j=1}^n W_{m+1,j} = \sum_{j=1}^n \beta^{M_j} \geq \beta^{M_i}$$

Note that M is number of mistakes, while m is number of running time. It is clear that $W_1 = n$, solving for M gives us:

$$M \leq \frac{\ln 1/\beta}{\ln 2/(1+\beta)} M_i + \frac{1}{\ln 2/(1+\beta)} \ln n$$

Setting $\beta = 1/e$ yields the result, completing the proof. \square

Remark 30. (Notion of Regret) We would like to obtain regret bound for arbitrary loss function $L : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty]$. Making our notation of regret more precise:

$$L_A(S) - \min_{i \in [n]} L_i(S) \leq o(m)$$

where $o(m)$ denotes some function that is sublinear in m that depends on other parameter:

$$\frac{L_A(S) - \min_{i \in [n]} L_i(S)}{m} \leq \frac{o(m)}{m}$$

Note that as $m \rightarrow \infty$:

$$\frac{L_A(S)}{m} \leq -\frac{\min_{i \in [n]} L_i(S)}{m}$$

The limit of the mean asymptotic loss is bounded by the mean of asymptotic loss of the best expert.

- The loss function $L : \{0, 1\} \times [0, 1] \rightarrow [0, \infty)$ the entropic loss given by:

$$L(y, \hat{y}) = y \ln \frac{y}{\hat{y}} + (1 - y) \ln \frac{1 - y}{1 - \hat{y}}$$

We can show that the regret with small constant with $\log(n)$

- Arbitrary loss function $L : \mathcal{Y} \times \hat{\mathcal{Y}} \rightarrow [0, B]$. The regret will be $\mathcal{O}(\sqrt{m \log n})$

Definition 5.5. (Simplex and Related Entropy) We define the simplex over probability distribution:

$$\Delta_n = \left\{ \mathbf{x} : [0, 1]^n : \sum_{i=1}^n x_i = 1 \right\}$$

We define the relative entropy $d : \Delta_n \times \Delta_n \rightarrow [0, \infty)$ as we have:

$$d(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^n u_i \ln \frac{u_i}{v_i}$$

we define the entropy loss is given as $L_{\text{en}}(y, \hat{y}) = d((y, 1 - y), (\hat{y}, 1 - \hat{y}))$

Definition 5.6. (Weighted Average) We will consider a projection in $[0, 1]$ rather than $\{0, 1\}$, and we will predict with weighted average. One weight per expert as we have $w_{t,i} = \beta^{L_{t,i}} = \exp(-\eta L_{t,i})$ where $L_{t,i}$ is the cumulative loss of E_i before the trial t , while η is the learning rate. The master algorithm predicts with the weighted average:

$$v_{t,i} = \frac{w_{t,i}}{\sum_{i=1}^n w_{t,i}} \quad \hat{y}_t = \sum_{i=1}^n v_{t,i} x_{t,i} = \mathbf{v}_t^T \mathbf{x}_t$$

where the $x_{t,i}$ is the prediction of E_i at trial t . We start with the initialize weight $\mathbf{v}_1 = \mathbf{w}_1 = (1/n, \dots, 1/n)$. This gives the pseudocode:

Algorithm 3 Weighted Average

- 1: **Input:** Input $\mathbf{v}_1 = \mathbf{w}_1 = (1/n, \dots, 1/n)$ with $L_{\text{WA}} = 0$ and $\mathbf{L} = \mathbf{0}$
- 2: **for** $i = 1, 2, \dots, m$ **do**
- 3: Receives instance $\mathbf{x}_t \in [0, 1]^n$
- 4: Predict $\hat{y}_t = \mathbf{v}_t^T \mathbf{x}_t$
- 5: Receives label $y_t \in [0, 1]$
- 6: Incur Loss $L_{\text{WA}} = L_{\text{WA}} + L(y_t, \hat{y}_t)$ and $L_i = L_i + L(y_t, x_{t,i})$ for $i \in [n]$
- 7: Update Weight, for $i \in [n]$:

$$v_{t+1,i} = \frac{v_{t,i} \exp(-\eta L(y_t, x_{t,i}))}{\sum_{j=1}^n v_{t,j} \exp(-\eta L(y_t, x_{t,j}))}$$

8: **end for**

Theorem 5.2. For sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in [0, 1]^n \times [0, 1]$. The regret of the weighted average WA algorithm is:

$$L_{\text{WA}}(S) - \min_i L_i(S) \leq 1/\eta \ln(n)$$

with square and entropic loss for $\eta = 1/2$ and $\eta = 1$ respectively.

Proof. We will proof the progress vs regret first, for all $\mathbf{u} \in \Delta_n$. Let's start with the assumption that $y_t = 1$

and by the error $L_{\text{en}}(1, x) = -\ln x$, we have

$$\begin{aligned}
d(\mathbf{u}, \mathbf{v}_t) - d(\mathbf{u}, \mathbf{v}_{t+1}) &= \sum_{i=1}^n u_i \ln \left(\frac{v_{t+1,i}}{v_{t,i}} \right) \\
&= \sum_{i=1}^n u_i \ln \frac{\frac{v_{t,i} \exp(-L_{\text{en}}(1, x_{t,i}))}{\sum_{j=1}^n v_{t,j} \exp(-L_{\text{en}}(1, x_{t,j}))}}{v_{t,i}} \\
&= \sum_{i=1}^n u_i \ln \frac{v_{t,i} x_{t,i}}{\sum_{j=1}^n v_{t,j} x_{t,j}} \\
&= \sum_{i=1}^n u_i \ln \frac{x_{t,i}}{\hat{y}_t} = \left(\sum_{i=1}^n u_i \ln x_{t,i} \right) - \ln(\hat{y}_t) \\
&= L_{\text{en}}(y_t, \hat{y}_t) - \sum_{i=1}^n u_i L_{\text{en}}(y_t, x_{t,i})
\end{aligned}$$

This also works by symmetry with the case $y = 0$, and so the claim:

$$d(\mathbf{u}, \mathbf{v}_t) - d(\mathbf{u}, \mathbf{v}_{t+1}) = L_{\text{en}}(y_t, \hat{y}_t) - \sum_{i=1}^n u_i L_{\text{en}}(y_t, x_{t,i})$$

is correct. Consider the telescoping sum, which we have:

$$\sum_{t=1}^m L_{\text{en}}(y_t, \hat{y}_t) - \sum_{t=1}^m \sum_{i=1}^n u_i L_{\text{en}}(y_t, x_{t,i}) = d(\mathbf{u}, \mathbf{v}_1) - d(\mathbf{u}, \mathbf{v}_{m+1})$$

Note that for any $\mathbf{u} \in \Delta_n$, especially the unit vector, which is shown to be an upper bound, and we can see that that $d(\mathbf{u}, \mathbf{v}_1) \leq \ln n$ and $-d(\mathbf{u}, \mathbf{v}_{m+1}) \leq 0$, and so we have proven the theorem, as:

$$\begin{aligned}
\sum_{t=1}^m L_{\text{en}}(y_t, \hat{y}_t) - \sum_{t=1}^m \sum_{i=1}^n u_i L_{\text{en}}(y_t, x_{t,i}) &= \sum_{t=1}^m L_{\text{en}}(y_t, \hat{y}_t) - \min_i L_i(S) \\
&\leq d(\mathbf{u}, \mathbf{v}_1) - d(\mathbf{u}, \mathbf{v}_{m+1}) \leq \ln(n) - 0
\end{aligned}$$

Note that we can have a unit vector u_i that is the correct expert. □

Definition 5.7. (Allocation Setting) On each trial, the learner plays an allocation $\mathbf{v}_t \in \Delta_t$, the nature returns the loss vector \mathbf{l}_t for example of the loss of expert i is $l_{t,i}$. There are 2 models for the learner:

- We can consider the incur loss directly: $L_A(t) = \mathbf{v}_t^T \mathbf{l}_t$
- The learner randomly select $\hat{y}_t \in [n]$ according to discrete distribution over $[n]$ with probability $v_{t,i}$ for each action, thus we have:

$$\mathbb{E}[L_A(t)] = \mathbb{E}[L_{t,\hat{y}}] = \mathbf{v}_t^T \mathbf{l}_t$$

The mechanism generating the loss vector \mathbf{l}_t must be obvious to the learner's selection \hat{y} until $t + 1$

This setting can simulate the setting where we rescue side-information \mathbf{x}_t and have a fixed loss function.

Theorem 5.3. (Hedge Theorem) For all sequence of loss vector $S = \mathbf{l}_1, \dots, \mathbf{l}_m \in [0, 1]^n$. The regret of the weighted average algorithm with $\eta = \sqrt{2m \ln n}$ is equal to:

$$\mathbb{E}[L_{\text{WA}}(S)] - \min_i L_i(S) \leq \sqrt{2m \ln n}$$

Proof. Given any $\mathbf{u} \in \Delta_n$. Letting $Z_t = \sum_{i=1}^n v_{t,i} \exp(-\eta l_{t,i})$, we observe that:

$$\begin{aligned}
d(\mathbf{u}, \mathbf{v}_t) - d(\mathbf{u}, \mathbf{v}_{t+1}) &= \sum_{i=1}^n u_i \ln \frac{v_{t+1,i}}{v_{t,i}} \\
&= -\eta \sum_{i=1}^n u_i l_{t,i} - \sum_{i=1}^n u_{t,i} \ln Z_t \\
&= -\eta \mathbf{u}^T \mathbf{l}_t - \ln \sum_{i=1}^n v_{t,i} \exp(-\eta l_{t,i}) \\
&\geq -\eta \mathbf{u}^T \mathbf{l}_t - \ln \sum_{i=1}^n v_{t,i} \exp\left(-\eta l_{t,i} + \frac{1}{2} \eta^2 l_{t,i}^2\right) \\
&= -\eta \mathbf{u}^T \mathbf{l}_t - \ln \left(1 - \eta \mathbf{v}_t^T \mathbf{l}_t + \frac{1}{2} \sum_{i=1}^n v_{t,i} l_{t,i}^2\right) \\
&\geq \eta(\mathbf{v}_t^T \mathbf{l}_t - \mathbf{u}^T \mathbf{l}_t) - \frac{1}{2} \eta^2 \sum_{i=1}^n v_{t,i} l_{t,i}^2
\end{aligned}$$

The first inequality uses $\exp(-x) \geq 1 - x + x^2/2$ and the second inequality uses $\ln(1+x) \leq x$. Now, let's consider the telescoping sum:

$$\begin{aligned}
\sum_{t=1}^m (\mathbf{v}_t^T \mathbf{l}_t - \mathbf{u}^T \mathbf{l}_t) &\leq \frac{1}{\eta} (d(\mathbf{u}, \mathbf{v}_1) - d(\mathbf{u}, \mathbf{v}_{m+1})) + \frac{\eta}{2} \sum_{t=1}^m \sum_{i=1}^n v_{t,i} l_{t,i}^2 \\
&\leq \frac{\ln n}{\eta} + \frac{\eta}{2} \sum_{t=1}^m \sum_{i=1}^n v_{t,i} l_{t,i}^2
\end{aligned}$$

This holds for all $\mathbf{u} \in \Delta_n$, it holds for a unit vector and we have the upper bound by noting that we have. $d(\mathbf{u}, \mathbf{v}_1) \leq \ln n$ and $-d(\mathbf{u}, \mathbf{v}_{m+1}) \leq 0$ and we have:

$$\sum_{t=1}^m \sum_{i=1}^n v_{t,i} l_{t,i}^2 \leq m$$

We can set $\eta = \sqrt{2 \ln n / m}$, as we have proven the theorem, as we can set $\eta = \sqrt{2 \ln n / m}$, which we have prove the theorem. \square

5.3 Online Learning of Linear Classifier

Definition 5.8. (Problem) We have the sequence of data: $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$ and the total loss is given by $L_A(S)$. The regret is defined as:

$$L_A(S) - \inf_{\mathbf{u} \in \mathcal{U}} \text{Loss}_{\mathbf{u}}(S)$$

where \mathcal{U} is a set of linear threshold function, as we will focus on the case where there exists $\mathbf{u} \in \mathcal{U}$ such that $\text{Loss}_{\mathbf{u}}(S) = 0$, which is a reliable case.

Definition 5.9. (Linear Threshold) The linear threshold $f_{\mathbf{u},b} : \mathbb{R}^n \rightarrow \{-1, 1\}$ function is:

$$f_{\mathbf{u},b}(\mathbf{x}) = \text{sgn}(\mathbf{u}^T \mathbf{x} + b)$$

The separating by hyperplane. The comparison class of all linear threshold function:

$$\mathcal{U}_{\text{it}} = \{f_{\mathbf{u},b} : \mathbf{u} \in \mathbb{R}^n, b \in \mathbb{R}\}$$

Remark 31. (Assumption) Data is linear separatable by some margin γ . Hence there exists a linear hyperplane with normal vector \mathbf{v} such that: $\|\mathbf{v}\| = 1$ and for all (\mathbf{x}_t, y_t) , which we have $y_t \in \{-1, 1\}$, and $\|\mathbf{x}_t\| \leq R$ and $y_t(\mathbf{x}_t^T \mathbf{v}) \geq \gamma$

Definition 5.10. (Perceptron Learning Algorithm) We consider the following learning algorithm

Algorithm 4 Perceptron Learning Algorithm

```

1: Initialize:  $\mathbf{w}_1 = \mathbf{0}$  and  $M_1 = 0$ 
2: for  $i = 1, 2, \dots, m$  do
3:   Receives Pattern  $\mathbf{x}_t \in \mathbb{R}^n$ 
4:   Predict  $\hat{y}_t = \text{sgn}(\mathbf{w}_t^T \mathbf{x}_t)$ 
5:   Receives Label  $y_t$ 
6:   if Mistake  $y_t \hat{y}_t \leq 0$  then
7:     Update  $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$ 
8:      $M_{t+1} = M_t + 1$ 
9:   else
10:     $\mathbf{w}_{t+1} = \mathbf{w}_t$  and  $M_{t+1} = M_t$ 
11:   end if
12: end for

```

Lemma 5.1. If $(\mathbf{w}_t^T \mathbf{x}_t)y_t < 0$ then $\|\mathbf{w}_{t+1}\|^2 \leq \|\mathbf{w}_t\|^2 + \|\mathbf{x}_t\|^2$

Proof. We have the following inequality:

$$\begin{aligned}
\|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t + y_t \mathbf{x}_t\|^2 \\
&= \|\mathbf{w}_t\|^2 + 2y_t(\mathbf{w}_t^T \mathbf{x}_t) + \|\mathbf{x}_t\|^2 \\
&\leq \|\mathbf{w}_t\|^2 + \|\mathbf{x}_t\|^2
\end{aligned}$$

□

Lemma 5.2. $\|\mathbf{w}_t\|^2 \leq M_t R^2$

Proof. Using induction, as we have:

- Base: $M_1 = 0$ and $\|\mathbf{w}_1\|^2 = 0$
- Induction step, when we have a mistake on the trial t as we have:

$$\|\mathbf{w}_{t+1}\|^2 \leq \|\mathbf{w}_t\|^2 + \|\mathbf{x}_t\|^2 \leq \|\mathbf{w}_t\|^2 + R^2 \leq (M_t + 1)R^2$$

If there is no mistake, then the outcome is trivial as we have $\mathbf{w}_{t+1} = \mathbf{w}_t$ and $M_{t+1} = M_t$

□

Lemma 5.3. $\|\mathbf{w}_t\|^2 \geq M_t \gamma$

Proof. Observe that $\|\mathbf{w}_t\| \geq \mathbf{w}_t^T \mathbf{v}$ because $\|\mathbf{v}\| = 1$ (via Cauchy-Schwarz). The prove of the lower bound $\mathbf{w}_t^T \mathbf{v}$ using the induction over t :

- Induction hypothesis $\mathbf{w}_t^T \mathbf{v} \geq M_t \gamma$
- Base $t = 1$, we have $\mathbf{w}_1^T \mathbf{v} = 0$
- Induction step: Assume for t and prove for $t + 1$, if there is a mistake as we have:

$$\begin{aligned}
\mathbf{w}_{t+1}^T \mathbf{v} &= (\mathbf{w}_t + \mathbf{x}_t y_t)^T \mathbf{v} \\
&= \mathbf{w}_t^T \mathbf{v} + y_t \mathbf{x}_t^T \mathbf{v} \\
&\geq M_t \gamma + \gamma = (M_t + 1) \gamma
\end{aligned}$$

This works in the case of non-mistake.

□

Theorem 5.4. For all sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \{-1, +1\}$. This mistake of the PERCEPTRON algorithm is bounded by:

$$M \leq \left(\frac{R}{\gamma}\right)^2$$

with $R = \max_t \|\mathbf{x}_t\|$. If there exists a vector \mathbf{v} with $\|\mathbf{v}\| = 1$ and constant γ such that $(\mathbf{v}^T \mathbf{x}_t)y_t \geq \gamma$

Proof. We use the bound on the norm of the weight $\|\mathbf{w}_t\|$ as we have:

$$(M\gamma)^2 \leq \|\mathbf{w}_{t+1}\|^2 \leq MR^2$$

and the inequality follows. □

Remark 32. It is convinence to express the bound in the form $M \leq R^2 \|\mathbf{u}\|^2$ where $\mathbf{u} = \mathbf{v}/\gamma$ then we have for all \mathbf{u} such that $(\mathbf{u}^T \mathbf{x}_t)y_t \geq 1$.

Remark 33. (Additional Problem) Suppose that \mathbf{w}_{m+1} doesn't necessary linearly separate S . How can we use the PERCEPTRON to define a vector \mathbf{w} and how long that would take ?

Remark 34. (Gradient Descent) Recalling the regularization approach to supervised learning as we have:

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{t=1}^m L(y_t, h(\mathbf{x}_t)) + \lambda \text{penalty}(h)$$

We consider the soft-margin SVM, which we have the following loss function:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}} \sum_{i=1}^m h_{\text{hi}}(y_i, \mathbf{w}^T \mathbf{x}_i + b) + \lambda \|\mathbf{w}\|^2$$

where $h_{\text{hi}}(y, \hat{y}) = \max(0, 1 - y\hat{y})$, which we can consider the following optimization problem:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^n} L_{\text{hi}}(y_i, \mathbf{w}^T \mathbf{x}_t) + \lambda \|\mathbf{w} - \mathbf{w}_t\|^2$$

Solving this problem, gives us:

$$\mathbf{w}_{t+1} = \begin{cases} \mathbf{w}_t & y_t(\mathbf{w}_t^T \mathbf{x}_t) > 1 \\ \mathbf{w}_t + \frac{y_t \mathbf{x}_t}{2\lambda} & y_t(\mathbf{w}_t^T \mathbf{x}_t) < 1 \end{cases}$$

Definition 5.11. (Online Gradient Descent) We consider the online gradient descent with hinge loss and $\|\cdot\|_2^2$ penalty, which we have the following pseudocode:

Algorithm 5 Online Gradient Descent

- 1: **Initialize:** $\mathbf{w}_1 = \mathbf{0}$ and $L_{\text{OGD}} = 0$
 - 2: Select $\eta \in (0, \infty)$
 - 3: **for** $i = 1, 2, \dots, m$ **do**
 - 4: Receives instance $\mathbf{x}_t \in \mathbb{R}^n$
 - 5: Predict $\hat{y}_t = \mathbf{w}_t^T \mathbf{x}_t$
 - 6: Receives Lable $y_t \in \{+1, -1\}$
 - 7: Incur Loss $L_{\text{OGD}} = L_{\text{OGD}} + L_{\text{hi}}(y_t, \hat{y}_t)$
 - 8: Update weight $\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbb{I}[y_t \hat{y}_t < 1] \eta y_t \mathbf{x}_t$
 - 9: **end for**
-

Theorem 5.5. Given $R = \max_t \|\mathbf{x}_t\|$ and $\|\mathbf{u}\| \leq U$. For the algorithm OGD with $\eta = U/(R\sqrt{m})$ as:

$$\sum_{t=1}^m L_{hi}(y_t, \hat{y}_t) - L_{hi}(y_t, \mathbf{u}^T \mathbf{x}_t) \leq \sqrt{U^2 R^2 m}$$

for any vector \mathbf{u} .

Proof. Using the convexity of the hinge loss (with respected to 2nd argument), which we have:

$$L_{hi}(y_t, \hat{y}_t) - L_{hi}(y_t, \mathbf{u}^T \mathbf{x}_t) \leq (\mathbf{w}_t - \mathbf{u})^T \mathbf{z}_t$$

where $\mathbf{z}_t = -y_t \mathbf{x}_t [y_t(\mathbf{w}_t^T \mathbf{x}_t) < 1] \in \partial_{\mathbf{w}} h_{hi}(y_t, \mathbf{w}_t^T \mathbf{x}_t)$. For the update, we have:

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 &= \|\mathbf{w}_t - \eta \mathbf{z}_t - \mathbf{u}\|^2 \\ &= \|\mathbf{w}_t - \mathbf{u}\|^2 - 2\eta(\mathbf{w}_t - \mathbf{u})^T \mathbf{z}_t + \eta^2 \|\mathbf{z}_t\|^2 \end{aligned}$$

And so, we have the:

$$(\mathbf{w}_t - \mathbf{u})^T \mathbf{z}_t = \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 + \eta^2 \|\mathbf{z}_t\|^2 \right)$$

and so we have:

$$\begin{aligned} \sum_{t=1}^m (\mathbf{w}_t - \mathbf{u})^T \mathbf{z}_t &= \sum_{t=1}^m \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 + \eta^2 \|\mathbf{z}_t\|^2 \right) \\ &\leq \frac{1}{2\eta} \left(\|\mathbf{u}\|^2 + \eta^2 \sum_{t=1}^m \|\mathbf{z}_t\|^2 \right) \\ &= \frac{1}{2\eta} \|\mathbf{u}\|^2 + \frac{\eta}{2} \sum_{t=1}^m \|\mathbf{x}_t\|^2 \mathbb{I}[y_t(\mathbf{w}_t^T \mathbf{x}_t) < 1] \\ &\leq \frac{1}{2\eta} U^2 + \frac{\eta}{2} m R^2 = \sqrt{U^2 R^2 m} \end{aligned}$$

as we have $\eta = U/(R\sqrt{m})$ and using the result from the convex setting yields the result. \square

Remark 35. (Perceptron Bound) The perceptron bound can be arrived by using the analysis of the OGD above as we have:

- If we consider the hinge, we have:

$$\sum_{t=1}^m \mathbb{I}[y_t \neq \text{sgn}(\hat{y}_t)] - L_{hi}(y_t, \mathbf{u}^T \mathbf{x}_t) \geq \sqrt{U^2 R^2 m}$$

- Assuming that there is a linear classifier \mathbf{u} such that $y_t(\mathbf{u}^T \mathbf{x}_t) \geq 1$ for all $t = 1, \dots, m$ as we have:

$$\sum_{t=1}^m \mathbb{I}[y_t \neq \text{sgn}(\hat{y}_t)] \geq \sqrt{U^2 R^2 m}$$

- Make OGD conservative that we only update when $y_t \hat{y}_t \leq 0$ instead of $y_t \hat{y}_t \leq 1$ as we have the trial when mistake is made.
- With respect to the bound, we can ignore the trial, which the mistake is made, so we can take the value $m = M := \sum_{t=1}^m \mathbb{I}[y_t \neq \text{sgn}(\hat{y}_t)]$, which implies that:

$$M \leq \sqrt{U^2 R^2 M} \implies M \leq U^2 R^2$$

- We can set $\eta = 1$ as its number doesn't matter at all. Recall the update rule for the perceptron, when mistake is made, with learning rate η : $\mathbf{w}_{t+1} = \mathbf{w}_t + \eta y_t \mathbf{x}_t$, as we have:

$$\mathbf{w} = \sum_{m:\text{mistake}} \eta y_t \mathbf{x}_t$$

If $\eta > 0$, then we have $\eta \sum_{m:\text{mistake}} y_t \mathbf{x}_t$. Note that the prediction made by perceptron is based on the sign of the dot product, and so η doesn't take on any effect.

5.4 Disjunction Learning

Definition 5.12. (Boolean Function) The boolean function f may be represented as a map $f : \{0, 1\}^n \rightarrow \{0, 1\}$ where, we have the following:

- $x_1 \wedge x_2 = x_1 x_2$
- $x_1 \vee x_2 = \text{sign}(x_1 + x_2)$
- $\bar{x} = 1 - x$

Furthermore, we have the following additional definition for the boolean function:

- Single variable is called a literal.
- Term or Conjunction is an iterated “and” applied.
- Clause or Disjunction is an iterated “or” applied.
- Monotone disjunction or conjunction implies no negated literal.

Remark 36. (Naive Weighted-Majority) The goal is to predict as well as k -literal (monotone) disjunction (over n variables). We can consider the use of weighted majority as each experts are disjunction of various variable and size. So, we have $\binom{n}{k}$ total expert and weights. This gives us the following bound:

$$\text{Mistake} \leq 2.63M + 2.63k \ln \frac{nl}{k}$$

where M is the mistakes of best disjunction, while we use the inequality

$$\binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$$

It is clear that the time and space are exponent in k (run time). We need better algorithm.

Corollary 5.1. *With the feature map $\phi(\mathbf{x}) = (\mathbf{x}, 1)$, we use the perceptron to learn monotone disjunction:*

$$M \leq (4k + 1)(n + 1)$$

when k is the number of literal out of the n possible literal. And, so there exists a generic lower bound for rotation invariance algorithm (SVM and perceptron) where $M = \Omega(n)$

Proof. We use the following perceptron bound:

$$M \leq R^2 \|\mathbf{u}\|^2$$

for all \mathbf{u} such that $(\mathbf{u}^T \mathbf{x}_t) y_t \geq 1$. With $\mathbf{x} \in \{0, 1\}^n$, the the feature map $\phi(\mathbf{x}) = (\mathbf{x}, 1)$, claim the following that $\mathbf{u}^* \in \mathbb{R}^{n+1}$ separate with margin of 1, where:

$$u_i^* = \begin{cases} 2 & i \text{ is a literal} \\ 0 & i \text{ isn't a literal} \\ -1 & i \text{ is bias weight} \end{cases}$$

Such that, we have the following calculation:

- $(\mathbf{u}^*)^T \phi(\mathbf{x}) \geq 1$ as we have positive example $y_t = 1$
- $(\mathbf{u}^*)^T \phi(\mathbf{x}) = -1$ as we have negative example $y_t = -1$

Note that for some $\mathbf{x} \in \{0, 1\}^n$, then $\|\phi(\mathbf{x})\|^2 \leq n + 1$ and we have $\|\mathbf{u}^*\|^2 = 4k + 1$, thus we have $M \leq (4k + 1)(n + 1)$ as required. \square

Definition 5.13. (Winnow Algorithm) We define the winnow algorithm to be

Algorithm 6 Winnow Algorithm

1: **Input:** $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m) \in \{0, 1\}^n \times \{0, 1\}$

2: **Initialize:** $\mathbf{w}_1 = \mathbf{1}$

3: Select $\eta \in (0, \infty)$

4: **for** $i = 1, 2, \dots, m$ **do**

5: Receives instances $\mathbf{x}_t \in \{0, 1\}^n$

6: Predict the value:

$$\hat{y}_t = \begin{cases} 0 & \mathbf{w}_t^T \mathbf{x}_t < n \\ 1 & \mathbf{w}_t^T \mathbf{x}_t \geq n \end{cases}$$

7: Receives the label $y_t \in \{0, 1\}$

8: **if** Mistake $\hat{y}_t \neq y_t$ **then**

9: Update the value:

$$w_{t+1,i} = w_{t,i} 2^{(y_t - \hat{y}_t)x_{t,i}} \quad \text{for } i \in [n]$$

10: **end if**

11: **end for**

Theorem 5.6. *The mistake of winnow is bounded by:*

$$M \leq 3k(\log n + 1) + 2$$

Proof. Let's consider 2 scenarios as we consider the bound on mistake:

- On a mistake, at least one element weight is doubled and the relevant weight never decreases.
- Once the weight $w_{t,i} \geq n$, it will no longer change (it will saturate to n) and so the mistake is:

$$M_p \leq k(\log n + 1)$$

where M_p is the bound on the positive example $y_t = 1$

Let $W_t = \sum_{i=1}^n w_{t,i}$. We can see that $W_1 = n$, and so:

- On the positive mistake ($y_t = 1$) we have $W_{t+1} \leq W_t + n$ (as we can only double it)
- On the negative mistake ($y_t = 0$) we have $W_{t+1} \leq W_t - n/2$ as we can only half the number of weights.

Consider the progression of weights, we can see that:

$$0 \leq W_{m+1} \leq W_1 + M_p n - M_f n/2 = n + M_p n - M_f n/2$$

Thus, we have $M_f \leq 2k(\log n + 1) + 2$, where M_p is the bound on the positive example $y_t = 0$. Combining them and we have:

$$M \leq M_p + M_f \leq 2 + 3k(\log n + 1)$$

\square

Remark 37. There are several observation that we have to make:

- WINNOW is an improvement over PERCEPTRON in terms of dimension m in the mistake bound.
- The bound for linear threshold learning for the WINNOW is incompatible as the algorithm prefer sparse hypothesis.

Theorem 5.7. *Given m , let t drawn uniformly at random from $\{1, \dots, m\}$. Let S be set of t examples sampled from p . Let (\mathbf{x}', y') be addiitonal example sample from P , then:*

$$\mathbb{P}(A_S(\mathbf{x}') \neq y') \leq \frac{B}{m}$$

with respected to be drawing of t, S and (\mathbf{x}', y) , where the mistake bound for A is B .

Proof. There are no more than B trials with mistake, therefore, since t is drawn uniformly from $\{1, \dots, m\}$ there is no more than B/m probability of hitting trial with a mistake. \square

Definition 5.14. (Disjunctive Normal Form) DNF is a disjunction of terms, for example:

$$x_1 x_4 x_7 \vee x_1 \bar{x}_2 \vee x_2 x_5$$

All boolean function may be represented as DNF.

Remark 38. DNF corresponds to simple boolean network with a single layer as such they may learn by a neural network with single hidden layer.

Definition 5.15. (ANOVA Kernel) We consider the feature map to be:

$$\mathbf{x} = \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix} \quad \Phi(\mathbf{x}) = \begin{matrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \\ x_1 x_2 \\ \vdots \\ x_1 x_2 \dots x_n \end{matrix}$$

There are 2^n features. The k -terms DNF in input space is k -literal in feature space:

$$\Phi(\mathbf{x})\Phi(\mathbf{y}) = \prod_{i=1}^n (1 + x_i y_i) = k_{\text{anova}}(\mathbf{x}, \mathbf{y})$$

Please note that it also represent a disjunction normal form.

Remark 39. (Perceptron for K-term DNF) The weight of the perceptron, we have the weight to be:

$$\mathbf{w}_t = \sum_{q \in \text{mistakes}} \alpha_q \Phi(\mathbf{x}_q)$$

And performing a predicting gives us:

$$\mathbf{w}_t^T \Phi(\mathbf{x}_t) = \left(\sum_{q \in \text{mistakes}} \alpha_q \Phi(\mathbf{x}_q) \right) \Phi(\mathbf{x}_t) = \sum_{q \in \text{mistakes}} \alpha_q \mathbf{k}(\mathbf{x}_q, \mathbf{x}_t)$$

The prediction time complexity is $\mathcal{O}(n \cdot \#\text{mistakes}) \leq \mathcal{O}(nm)$. The mistake bound is $\mathcal{O}(k2^n)$

Remark 40. The winnow weight is given as:

$$w_{t,i} = \exp \left(-\eta \sum_{q \in \text{mistake}} \alpha_q [\Phi(\mathbf{x}_q)]_i \right)$$

We have the log of weights that is linear combination of the past examples. We have the mistake bound to be $\mathcal{O}(k \ln 2^n) = \mathcal{O}(kn)$ with the prediction time to be $\Omega(2^n \# \text{mistake})$, as there is no obvious fast way to compute $\mathbf{w}_t^T \Phi(\mathbf{x}_t)$ for now.

6 Online Learning 2: Bandits

Definition 6.1. (Partial Feedback Protocol) We consider the following setting:

Algorithm 7 Partial Feedback Control

```

1: for  $i = 1, 2, \dots, T$  do
2:   Predict  $\hat{y}_t \in [n]$ 
3:   Observe loss of prediction  $l_{t,\hat{y}_t} \in [0, 1]$ 
4: end for
```

We have the following goal:

$$\sum_{t=1}^m l_{t,\hat{y}_t} - \min_{i \in [n]} \sum_{t=1}^m l_{t,i} \leq o(m)$$

This is the same as the regret. Please note that we didn't get to see all loss function that is induced by the prediction.

Definition 6.2. (Unbiased Estimation) An estimator $\hat{\theta}$ estimate a parameter θ of a distribution from a sample is unbiased if we can show that $\mathbb{E}[\hat{\theta}] = \theta$.

Example 6.1. Suppose X_1, \dots, X_n are iid random variable for a distribution with mean μ , then:

$$\hat{\theta} = \frac{1}{n}(X_1 + \dots + X_n)$$

is an unbiased estimate of μ

Example 6.2. Suppose X is a random variable with the discrete uniform distribution over $\{1, \dots, n\}$. Suppose n is unknown and we wish to estimate it.

- The estimate $\hat{\theta}_1 = X$ is the maximum likelihood estimator, since $\mathcal{L}(\theta, X = x) = 1/\theta$ is maximized when $\theta = x$. Then we have:

$$\mathbb{E}[\hat{\theta}_1; \theta = n] = \sum_{x=1}^n \frac{x}{n} = \frac{n+1}{2}$$

- And so, $\hat{\theta}_2 = 2x - 1$ is unbiased estimator, which is:

$$\mathbb{E}[\hat{\theta}_2; \theta = n] = \sum_{x=1}^n \frac{1}{n}(2x - 1) = 2 \sum_{x=1}^n \frac{1}{n}(2x - 1) = 2 \sum_{x=1}^n \frac{1}{n}x - \sum_{x=1}^n \frac{1}{n} = n$$

Remark 41. (Assumption and Estimation) Suppose, we have a distribution D_i over $[0, 1]$ for each $i \in [n]$ arms. For each arm i , we use iid sample $l_{t,i}$ for D_i . Suppose, we play i on trials $S_{t,i} \subseteq [t]$, then:

$$\hat{\mu}_{t,i} = \sum_{t \in S_i} \frac{l_{t,i}}{|S_i|}$$

This is unbiased estimator of μ_i . Now, we can consider the usage as we have:

- We can use a concentration inequality that allows us to quantitatively estimate the likelihood to estimate differently for the parameter.
- Using the observation, the algorithm UCB balances exploration and exploitation to obtain good regret bounds for this method.
- Suppose the underlying D_i is changing over time (being $D_{t,i}$):

$$\mu_{t,i} = \frac{\sum_{j=1}^t \mathbb{E}[l_{j,t}]}{t}$$

where $S_i = [t]$. However, if we only have $S_{t,i} = [t]$, then we have no information about the other arms.

- We need to have simultaneous unbiased estimate for all arms S

Definition 6.3. (Importance Weighting) We have the following series of observation:

- Suppose X is a random variable over \mathbb{R} with a mean μ . By definition, $\mathbb{E}[X] = \mu$ and $\hat{\theta}_1 = X$ is an unbiased estimator of the mean.
- Consider the biased coin Z_p with outcome 1 with probability p . Suppose, we have the estimator $\hat{\theta}_0$ setting to equal to X/p if $Z_p = 1$.
- Its expectation is equal to:

$$\mathbb{E}[\hat{\theta}_0] = \mathbb{P}(Z_p = 1)(X/p) + 0\mathbb{P}(Z_p = 0) = (p)(X/p) + (1-p)(0) = X$$

This is unbiased.

Definition 6.4. (Hallucinated Loss Vector) We generalize this to obtain an unbiased estimator of l_t in the bandit setting. Given $\mathbf{v}_t \in \Delta_n$ by the relation that $\hat{y}_t \sim \mathbf{v}_t$. The unbiased estimator l_t^h or \mathbf{k}_t with respect to \mathbf{v}_t is given as:

$$\left(l_{t,i}^h = \frac{l_{t,i}}{v_{t,i}} \mathbb{I}[i = \hat{y}_t] \right)_{i \in [n]}$$

Remark 42. (Expectation of Hallucinated Loss Vector) Observed that l_t^h is unbiased for all $i \in [n]$ since we have:

$$\mathbb{E}_{\hat{y}_t, \mathbf{v}_t} [l_{t,i}^h] = \sum_{j=1}^n v_{t,j} \frac{l_{t,i}}{v_{t,i}} \mathbb{I}[i = j] = l_{t,i}$$

We have unbiased estimator for all arms by only observing the single arm. We can apply the hedge to l_t^h requires bounded loss vector. We can use more careful analysis of the hedge.

Definition 6.5. (EXP3) Exponential-Weight algorithm for Exploration and Exploitation is given by:

Algorithm 8 EXP3

- 1: **Initialize:** $\eta \in (0, \infty)$
- 2: Set $\mathbf{v}_1 = (1/n, \dots, 1/n)$
- 3: **for** $i = 1, 2, \dots, T$ **do**
- 4: Sample $\hat{y}_t \sim \mathbf{v}_t$
- 5: Observe Loss $l_{t,\hat{y}} \in [0, 1]$
- 6: Construct Hallucinated Loss vector:

$$l_t^h = \left(l_{t,i}^h = \frac{l_{t,i}}{v_{t,i}} \mathbb{I}[i = \hat{y}_t] \right)_{i \in [n]}$$

- 7: Perform the update, for $i \in [n]$ and $Z_t = \sum_{i=1}^n v_{t,i} \exp(-\eta l_{t,i}^h)$:

$$v_{t+1,i} = v_{t,i} \exp(-\eta l_{t,i}^h) / Z_t$$

- 8: **end for**
-

Lemma 6.1. For any sequence of loss vector $\mathbf{l}_1, \dots, \mathbf{l}_m \in [0, 1]^n$, we have the following loss bound:

$$\sum_{t=1}^m \mathbf{v}_t^T \mathbf{l}_t^h - \sum_{t=1}^m \mathbf{u}^T \mathbf{l}_t^h \leq \frac{\ln n}{\eta} + \frac{\eta}{2} \sum_{t=1}^m \sum_{i=1}^n v_{t,i} (l_{t,i}^h)^2$$

For all $\mathbf{u} \in \Delta_n$

Proof. The lemma follows from the fact that EXP3 is just Hedge with \mathbf{l}_t weighted to be \mathbf{l}_t^h and the Hedge inequality is proven before. \square

Remark 43. We can show the property of EXP3, where we consider that: we need to perform and so we may replace hallucination losses \mathbf{l}_t^h with time loss \mathbf{l} :

- We can model some of the randomness as we use the adversarial loss $\mathbf{l}_1, \dots, \mathbf{l}_m$.
- We have to bound the term $\sum_{t=1}^m \sum_{i=1}^n v_{t,i} (l_{t,i}^h)^2$ and tune η

Definition 6.6. (Deterministic Adversarial Model) We will to set $\mathbf{l}_1, \dots, \mathbf{l}_m$ before running the algorithm. The adversary is assumed to be complete given the prior knowledge, and:

- The limitation of near omniscient adversary is that it is non-adaptive.
- It many simulate the stochastic model by repeatedly sample the $\mathcal{D}_1, \dots, \mathcal{D}_m$ in advance.

Theorem 6.1. For any sequence of loss vector $S = \mathbf{l}_1, \dots, \mathbf{l}_m \in [0, 1]^n$, the regret for EXP3 with $\eta = \sqrt{2 \ln n / mn}$ is:

$$\mathbb{E}[L_A(S)] - \min_i L_i \leq \sqrt{2mn \ln n}$$

where $L_A(S) = \sum_{t=1}^m l_{t, \hat{y}_t}$ and $L_i = \sum_{t=1}^m l_{t,i}$

Proof. Observe that the only source of randomness are the sample $\hat{y}_t \sim \mathbf{v}_t$. As previously argue, note that $\mathbb{E}[l_{t,i}^h] = l_{t,i}$, and we have:

$$\mathbb{E}[\mathbf{v}_t^T \mathbf{l}_t^h] = \sum_{i=1}^n \mathbb{E}[v_{t,i} l_{t,i}^h] = \sum_{i=1}^n v_{t,i} \mathbb{E}[l_{t,i}^h] = \sum_{i=1}^n v_{t,i} l_{t,i} = \mathbb{E}[l_{t, \hat{y}_t}]$$

Similarly, we have:

$$\mathbb{E}[(l_{t,i}^h)^2] = \sum_{j=1}^n v_{t,j} \left(\frac{l_{t,i}}{v_{t,i}} \right)^2 \mathbb{I}[i=j] = v_{t,i} \left(\frac{l_{t,i}}{v_{t,i}} \right)^2 = \frac{l_{t,i}^2}{v_{t,i}}$$

This implies that:

$$\mathbb{E} \left[\sum_{i=1}^n v_{t,i} (l_{t,i}^h)^2 \right] = \sum_{i=1}^n v_{t,i} \frac{l_{t,i}^2}{v_{t,i}} = \sum_{i=1}^n l_{t,i}^2 \leq n$$

Taking the expectation over the Hedge terms, and we have for $\mathbf{u} \in \Delta_n$:

$$\mathbb{E} \left[\sum_{t=1}^m \mathbf{v}_t^T \mathbf{l}_t^h - \sum_{t=1}^m \mathbf{u}^T \mathbf{l}_t^h \right] \leq \mathbb{E} \left[\frac{\ln n}{\eta} + \frac{\eta}{2} \sum_{t=1}^m \sum_{i=1}^n v_{t,i} (l_{t,i}^h)^2 \right]$$

And, so we have using the fact that: $\mathbb{E}[l_{t,i}^h] = l_{t,i}$, and the previous result with \mathbf{u} being a coordinate vector, we have:

$$\mathbb{E} \left[\sum_{t=1}^m \mathbf{v}_t^T \mathbf{l}_t^h \right] - \min_i \mathbb{E} \left[\sum_{t=1}^m l_{t,i}^h \right] \leq \frac{\ln n}{\eta} + \frac{\eta}{2} \mathbb{E} \left[\sum_{t=1}^m \sum_{i=1}^n v_{t,i} (l_{t,i}^h)^2 \right]$$

And, so we have:

$$\mathbb{E}[L_A(S)] - \min_i L_i(S) \leq \ln \frac{n}{\eta} + \frac{\eta}{2} mn$$

Substitute the $\eta = \sqrt{2 \ln n / mn}$ to prove this theorem. \square

7 Learning Theory

7.1 Introduction

Definition 7.1. (Distribution over Subset) If \mathcal{D} is a distribution over \mathcal{Z} then if $A \subseteq \mathcal{Z}$ then $\mathcal{D}(A)$ denotes the probability that if z is drawn from \mathcal{D} that $z \in A$

Definition 7.2. (Expected Error) Data is sampled iid from a distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$ with $\mathcal{Y} = \{0, 1\}$. The expected error of function $h : \mathcal{X} \rightarrow \mathcal{Y}$ is:

$$L_{\mathcal{D}}(h) = \mathcal{D}(\{x, y\} : h(x) \neq y) = \mathbb{P}_{(x,y) \sim \mathcal{D}}[h(x) \neq y] = \int \mathbb{I}[h(x) \neq y] \, d\mathcal{P}(x, y)$$

where we denote $L_{\mathcal{D}}(h) = \mathcal{E}(h)$

Definition 7.3. (Empirical Error) The empirical error of h given the dataset $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ is denoted as:

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq y_i]$$

Or, we denote it as $\mathcal{E}_{\text{emp}}(S, h)$.

Theorem 7.1. (Hoeffding's Inequality) Let Z_1, Z_2, \dots, Z_m be iid bernoulli random variable, when for all i , we have $\mathbb{P}(Z_i = 1) = p$ and let $\bar{Z} = 1/m \sum_{i=1}^m Z_i$, then for any $\varepsilon > 0$ as we have:

$$\mathbb{P}(\bar{Z} > p + \varepsilon) \leq \exp(-2m\varepsilon^2) \quad \mathbb{P}(\bar{Z} < p - \varepsilon) \leq \exp(-2m\varepsilon^2)$$

Theorem 7.2. Select a function h then for any $\delta \in (0, 1)$ with probability $1 - \delta$ over the random sample V of size m from \mathcal{D} , we have:

$$L_{\mathcal{D}}(h) \leq L_V(h) + \sqrt{\frac{\ln(1/\delta)}{2m}}$$

The generalization error of a function h may be bounded by the empirical error. We may select a predictor h on any set S , as we may bound it on the validation on separate set of data V .

Proof. Given a predictor h , we have the differences to be:

$$L_{\mathcal{D}}(h) - L_V(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}}[h(x) \neq y] - \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(x_i) \neq y_i]$$

we can define $Z_i = \mathbb{I}[h(x_i) \neq y_i]$. We can see that Z_1, \dots, Z_m are statistical independent. Then for all $\mathbb{P}[Z_i] = L_{\mathcal{D}}(h) = \mathbb{P}[h(x) \neq y]$. We apply the Hoeffding inequality, gives us:

$$\mathbb{P}[L_{\mathcal{D}}(h) - L_V(h) \geq \varepsilon] \leq \exp(2 - \varepsilon^2)$$

setting $\delta = \exp(-2\varepsilon^2 m)$, and solving this gives us the theorem. □

Remark 44. If we use the upper and lower bound m , the Hoeffding inequality would gives us:

$$|L_{\mathcal{D}}(h) - L_V(h)| \leq \sqrt{\frac{\ln(2/\delta)}{2m}}$$

This is a nice result, but there are some drawbacks to this bound:

- The validation bound gives a way to estimate of the confidence interval for the generalization error. The data V can't be used for training.
- Having small number of data, can we choose a model based on the expected error directly, without the training data ?

- The bound is about the predictor, while we need to analyze the prediction done by the machine learning algorithm.

Definition 7.4. (General Statistical Consider) Statistical model begin with an assumption that the data is generated by the underlying distribution \mathcal{D} not known to the learner. Assuming that we are given a training set that is generated iid from distribution \mathcal{D} :

$$S = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

Definition 7.5. (Empirical Risk Mimimization) Assuming we have a learning algorithm A that chooses a hypothesis function $A_{\mathcal{H}}(S)$ from a hypothesis space \mathcal{H} in response to the training set S . We study the ERM:

$$\text{ERM}_{\mathcal{H}}(S) = \arg \min_{h \in \mathcal{H}} L_S(h)$$

There are many possible empirical minimizer as we assume ERM to be an arbitrary one.

Remark 45. The traditional statistic h_S concentrated on analysing:

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S_m} [L_{\mathcal{D}}(A(S_m))]$$

where S_m denotes a training set of size m . For finite sample, the generalization $L_{\mathcal{D}}(A(S_m))$ has a distribution depending on the algorithm and function class and sample size:

- *Traditional Statistic:* concentrated on the mean of this distribution but this quantity is misleading for example in the case of low fold cross-validation.
- *Statistical Learning Theory:* analyze the tail of the distribution finding and the bound that holds in high probability.

Definition 7.6. (Reliability Assumption) Assume that there exists a function f^* so that for all $x \in \mathcal{X}$, we have $f^*(x) = y$ there exists a classifier that has zero error. We can now take \mathcal{D} to be only a distribution over \mathcal{X} only. We consider the following loss:

$$L_{\mathcal{D}, f^*}(h) = \mathbb{P}[h(x) \neq f^*(x)]$$

We can find the algorithm A so that $h = A(S)$ such that $L_{\mathcal{D}, f^*}(h) = 0$ is small.

Remark 46. (Reason for Approximation) We can't hope to find the function h such that $L_{\mathcal{D}, f^*}(h) = 0$. Let's consider the $\varepsilon \in (0, 1)$ that takes $\mathcal{X} = \{x_1, x_2\}$ where $\mathcal{D}(\{x_1\}) = 1 - \varepsilon$ and $\mathcal{D}(\{x_2\}) = \varepsilon$:

- The probability to not see x_2 at all among m iid example is $(1 - \varepsilon)^m \approx \exp(-\varepsilon m)$
- If $\varepsilon \ll 1/m$, we are unlikely to see x_2 at all. then we don't know its label.

So, we are only happy to see $L_{\mathcal{D}, f^*}(h) \leq \varepsilon$ when ε is user defined.

Remark 47. (Reason for Probability) The input is randomly generated (there is a small chance that we will see the same sample over and over again). No algorithm can generate $L_{\mathcal{D}, f^*}(h) \leq \varepsilon$ for sure, and so we allow the algorithm to fail with some probability $\delta \in (0, 1)$ that is user defined.

Definition 7.7. (PAC Learning) The learner doesn't know \mathcal{D} and f^* . It receives parameter ε and δ . Learner can ask for training data S contrary for $m(\varepsilon, \delta)$ examples. The learner should output a hypothesis h such that with at least probability $1 - \delta$, it holds that $L_{\mathcal{D}, f^*}(h) \leq \varepsilon$

Theorem 7.3. (No Free Lunch)

- Suppose $|\mathcal{X}| = \infty$. For any fixed $C \subset \mathcal{X}$ take \mathcal{D} to be uniform m distribution over C :
- If the number of training example is $m \leq |C|/2$, the learner has no knowledge of at least half of elements in C .

Fix $\delta \in (0, 1)$ and $\varepsilon < 1/2$. For any learner A and training set of size m , there exists \mathcal{D} and f^* such that with probability δ over the generation of a training data S of m examples, it holds that

$$L_{\mathcal{D}, f^*}(A(S)) \geq \varepsilon$$

Proof. Consider for contradiction, assuming that the class is learnable, consider $\varepsilon > 1/8$ and $\delta \leq 1/7$. With the definition of PAC learnable $m(\varepsilon, \delta) = m$:

- For the consistent case, with probability greater than $1 - \delta$, when A is applied to sample S of size m , generated iid \mathcal{D} , we have

$$L_{\mathcal{D}, f^*}(A(S)) \leq \varepsilon$$

- However, using the NFL theorem above, since $|\mathcal{X}| > 2m$, for every learning algorithm, there exists a \mathcal{D} such that with probability greater than $1/7 > \delta$, and $L_{\mathcal{D}, f^*}(A(S)) > 1/8 > \varepsilon$

This is a contradiction. \square

7.2 PAC of Finite Hypothesis Class

Lemma 7.1. For any 2 sets A and B , and a distribution \mathcal{D} we can show that:

$$\mathcal{D}(A \cup B) \leq \mathcal{D}(A) + \mathcal{D}(B)$$

Theorem 7.4. Fix ε, δ . If we have $m \geq \log(|\mathcal{H}|/\delta)/\varepsilon$, then for every \mathcal{D}, f^* with probability of at least $1 - \delta$ (with respect to randomly sample training set S of size m), we now have:

$$L_{\mathcal{D}, f^*}(ERM_{\mathcal{H}}(S)) \leq \varepsilon$$

This mean that we have $L_{\mathcal{D}, f^*}(ERM_{\mathcal{H}}(S)) \leq (\log|\mathcal{H}| + \log(1/\delta))/m$. The generalization error decrease linear in the number of samples and increase in logarithm in the size of hypothesis class.

Proof. Consider $S|_x = (x_1, \dots, x_m)$ be instances of training set. We will show that:

$$\mathcal{D}^m(\{S|_x : L_{\mathcal{D}, f^*}(ERM_{\mathcal{H}}(S)) > \varepsilon\}) \leq \delta$$

Let \mathcal{H}_B be a set of bound hypothesis as we have $\mathcal{H}_B = \{h \in \mathcal{H} : L_{\mathcal{D}, f^*}(h) > \varepsilon\}$ and let M be the set of misleading samples: $\{S|_x : \exists h \in \mathcal{H}_B, L_S(h) = 0\}$ Observe that:

$$\{S|_x : \exists h \in \mathcal{H}_B, L_S(h) = 0\} \subseteq M = \bigcup_{h \in \mathcal{H}_B} \{S|_x : L_S(h) = 0\}$$

Applying the union bound as we have the following union bound:

$$\begin{aligned} & \mathcal{D}^m(\{S|_x : \exists h \in \mathcal{H}_B, L_S(h) = 0\}) \\ & \leq \sum_{h \in \mathcal{H}_B} \mathcal{D}^m(\{S|_x : L_S(h) = 0\}) \\ & \leq |\mathcal{H}_B| \max_{h \in \mathcal{H}_B} \mathcal{D}^m(\{S|_x : L_S(h) = 0\}) \\ & < |\mathcal{H}_B|(1 - \varepsilon)^m \leq |\mathcal{H}| \exp(-\varepsilon m) \end{aligned}$$

Observe that $\mathcal{D}^m(\{S|_x : L_S(h) = 0\}) = (1 - L_{\mathcal{D}, f^*}(h))^m$ if $h \in \mathcal{H}_B$, then $L_{\mathcal{D}, f^*}(h) \geq \varepsilon$. This leads to the third inequality, while the last inequality, we have: $1 - \varepsilon \leq \exp(-\varepsilon)$ and $|\mathcal{H}_B| \leq |\mathcal{H}|$. Setting the rhs to $\leq \delta$ and we get the required inequality. \square

Definition 7.8. (PAC-Learnability) A hypothesis class \mathcal{H} is PAC-learnable if there exists a function $m_{\mathcal{H}} : (0, 1)^2 \rightarrow \mathbb{N}$ and has the property of that for every ε and $\delta \in (0, 1)$ and every distribution \mathcal{D} over \mathcal{X} and for every labeling function $f^* : \mathcal{X} \rightarrow \{0, 1\}$.

- Using the training algorithm $m \geq m_{\mathcal{H}}(\varepsilon, \delta)$ iid examples generated by \mathcal{D} and labeled by f^*
- The algorithm returns a hypothesis h such that with probability of at least $1 - \delta$, the loss is $L_{\mathcal{D}, f^*}(h) \leq \varepsilon$.
- We call $m_{\mathcal{H}}$ is the sample complexity of the training hypothesis \mathcal{H}

Remark 48. We are now interested in the infinite hypothesis space. What is the sample complexity of a given class ? Is there a generic algorithm that achieves the optimal sample complexity ?

Remark 49. (VC-Dimension: Motivation) Suppose, we have the training set: $S = (x_1, y_1), \dots, (x_m, y_m)$. We try to explain the label using a hypothesis from \mathcal{H} . We may get difference labels:

$$(x_1, y'_1), \dots, (x_m, y'_m)$$

We can try to explain the label using a hypothesis from \mathcal{H} . If this works for us, no matter the labels are then, no free-lunch theorem apply, as now we can't learn from $m/2$ example.

Definition 7.9. (VC-Dimension) Let $C = \{x_1, \dots, x_{|C|}\} \subset \mathcal{X}$. Let \mathcal{H}_C is the restriction of \mathcal{H} to C , then we have:

$$\mathcal{H}_C = \{h_C : h \in \mathcal{H}\} \quad \text{where} \quad h_C : C \rightarrow \{-1, 1\}$$

is such that $h_C(x_i) = h(x_i)$. For every $x_i \in C$, we can represent each h_C as the vector:

$$\mathcal{H}_C = \left\{ (h(x_1), \dots, h(x_{|C|})) \in \{-1, 1\}^{|C|} \right\}$$

and so we have $|\mathcal{H}_C| \leq 2^{|C|}$. We say that \mathcal{H} shatters C if $|\mathcal{H}_C| = 2^{|C|}$ where we have:

$$\text{VCDim}(\mathcal{H}) = \sup \{|C| : \mathcal{H} \text{ shatters } C\}$$

VC dimension is the maximum size of a set C such that \mathcal{H} gives no prior knowledge with respect to C .

Remark 50. To show that the VC dimension $\text{VCDim}(\mathcal{H}) = d$, we have to show that:

- There exists a set C of size d which is shattered by \mathcal{H}
- Every set C of size $d + 1$ isn't shattered by \mathcal{H}

Proposition 7.1. (VC-Dimension of Intervals) Interval where we have $\mathcal{H} = \mathbb{R}$ and

$$\mathcal{H} = \{h_{a,b} : a < b \in \mathbb{R}\}$$

where $h_{a,b}(x) = 1$ iff $x \in [a, b]$. Its VC-Dimension is 2.

Proposition 7.2. (Axis Aligned Rectangle) We have $\mathcal{X} = \mathbb{R}^2$ as we have the hypothesis set to be:

$$\mathcal{H} = \{h_{(a_1, a_2, b_1, b_2)} : a_1 < a_2 \text{ and } b_1 < b_2\}$$

where we have $h_{(a_1, a_2, b_1, b_2)}(\mathbf{x}_1, \mathbf{x}_2) = 1$ iff $\mathbf{x}_1 \in [a_1, a_2]$ and $\mathbf{x}_2 \in [b_1, b_2]$. We can show that $\text{VCDim}(\mathcal{H}) = 4$

Proof. We can find 4 points that can be shattered by \mathcal{H} , and so $\text{VCDim}(\mathcal{H}) \geq 4$. For any point $C \subseteq \mathbb{R}^2$ with 5 points with label $(1, 1, 1, 1, 0)$ where 0 is the point in the middle, we can't obtain any axis aligned rectangle, thus it can't be shattered C . Therefore, $\text{VCDim}(\mathcal{H}) = 4$ \square

Proposition 7.3. (Finite Class) The VC-Dimension of the finite \mathcal{H} is at most $\log_2(|\mathcal{H}|)$ as there can arbitrary gaps between $\text{VCDim}(\mathcal{H})$ and $\log_2(|\mathcal{H}|)$

Proof. Let \mathcal{H} be a finite class, for any set C that can be shattered, we have $2^{|C|} = |\mathcal{H}_C| \leq |\mathcal{H}|$, thus the upperbound of the VC dimension is $\log_2 |\mathcal{H}|$ \square

Theorem 7.5. (Radon) Any set \mathcal{X} of $d + 2$ data point \mathbb{R}^d can be partition into 2 sets \mathcal{X}_1 and \mathcal{X}_2 such that the convex hull of \mathcal{X}_1 and \mathcal{X}_2 intersect.

Proof. Let $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{d+2}\} \subset \mathbb{R}^d$ with the following linear equation:

$$\sum_{i=1}^{d+2} \alpha_i \mathbf{x}_i = 0 \quad \sum_{i=1}^{d+2} \alpha_i = 0$$

The number of unknown $d+2$ is larger than the number of equations $d+1$. This implies that the system admits non-zero solution $\beta_1, \dots, \beta_{d+2}$ since $\sum_{i=1}^{d+2} \beta_i = 0$ both:

$$\mathcal{J}_1 = \{i \in [d+2] : \beta_i > 0\} \quad \mathcal{J}_2 = \{i \in [d+2] : \beta_i \leq 0\}$$

This means that $\mathcal{X}_1 = \{\mathbf{x}_i : i \in \mathcal{J}_1\}$ and $\mathcal{X}_2 = \{\mathbf{x}_i : i \in \mathcal{J}_2\}$ form a partition. The last equation gives us:

$$\sum_{i \in \mathcal{J}_1} \beta_i = - \sum_{i \in \mathcal{J}_2} \beta_i$$

Let $\beta = \sum_{i \in \mathcal{J}_1} \beta_i$, then the first equation implies that:

$$\sum_{i \in \mathcal{J}_1} \frac{\beta_i}{\beta} \mathbf{x}_i = - \sum_{i \in \mathcal{J}_2} \frac{\beta_i}{\beta} \mathbf{x}_i$$

Please note that: $\sum_{i \in \mathcal{J}_1} \beta_i/\beta = -\sum_{i \in \mathcal{J}_2} \beta_i/\beta = 1$ and $\beta_i/\beta \geq 0$ for $i \in \mathcal{J}_1$ and $-\beta_j/\beta \geq 0$ for $i \in \mathcal{J}_2$. By the definition of the convex hull, this implies that $\sum_{i \in \mathcal{J}_1} \beta_i/\beta \mathbf{x}_i$ being both to convex hull \mathcal{X}_1 and \mathcal{X}_2 \square

Proposition 7.4. (Hyperplane) We have $\mathcal{X} = \mathbb{R}^n$ and the hypothesis class to be:

$$\mathcal{H} = \{y \mapsto \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle) : \mathbf{w} \in \mathbb{R}^n\}$$

Then, we have $\text{VCDim}(\mathcal{H}) = n + 1$

Proof. Starting with the lower bound, setting \mathbf{x}_0 to be the origin and setting \mathbf{x}_i for $i \in [d]$ as the whose i coordinate to be 1 and all the others are 0.

- Let $y_0, y_1, \dots, y_d \in \{-1, 1\}$ be an arbitrary set of label.
- Let \mathbf{w} be the vector whose i -th coordinate is y_i .

The classifier defined by the hyperplane of equation $\mathbf{w}^T \mathbf{x} + y_0/2 = 0$, shatters $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_d$ we can see that for any $i \in \{0, \dots, d\}$ as we have:

$$\text{sgn}\left(\mathbf{w}^T \mathbf{x}_i + \frac{y_0}{2}\right) = \text{sgn}\left(y_i + \frac{y_0}{2}\right) = y_i$$

For the upperbound, let \mathcal{X} be set of $d+2$ points. By Radon's theorem, it can be partition into 2 sets \mathcal{X}_1 and \mathcal{X}_2 such that the convex hull intersects. When the set of points \mathcal{X}_1 and \mathcal{X}_2 are separated by hyperplane, the convex hull also separated. However, it is a contradiction and so the VC dimension is proven. \square

Definition 7.10. (Inner Production Space) The space is the bounded sequence summable square:

$$l_2 = \left\{ \mathbf{x} \in \mathbb{R}^\infty : \sum_{i=1}^{\infty} x_i^2 < \infty \right\}$$

with the inner product to be $\{\mathbf{x}, \mathbf{x}'\} = \sum_{i=1}^{\infty} x_i x'_i$

Definition 7.11. (Large Margin Halfspaces) Given $\mathcal{X} \subset l_2$ and $\Lambda \in (0, \infty)$, which we define:

$$\mathcal{H}_{\mathcal{X}, \Lambda} = \{\mathbf{x} \mapsto \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle) : \mathbf{x} \in \mathcal{X}, \mathbf{w} \in l_2, \|\mathbf{w}\| \leq \Lambda, \langle \mathbf{w}, \mathbf{x} \rangle \geq 1\}$$

Observe that $1/\|\mathbf{w}\|$ is the margin.

Theorem 7.6. *We can show that for large margin halfspace:*

$$\text{VCDim}(\mathcal{H}_{\mathcal{X},\Lambda}) \leq \Lambda^2 \max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_{\mathcal{X}}^2$$

Theorem 7.7. (Fundamental Theorem of Statistical Learning) *Let \mathcal{H} be a hypothesis class of binary classifier. Then there are absolute constant C_1 and C_2 such that the sample complexity is given by:*

$$C_1 \frac{\text{VCDim}(\mathcal{H}) + \log(1/\delta)}{\varepsilon} \leq m_{\mathcal{H}}(\varepsilon, \delta) \leq C_2 \frac{\text{VCDim}(\mathcal{H}) \log(1/\varepsilon) + \log(1/\delta)}{\varepsilon}$$

This sample complexity is achieved by ERM learning rule.

7.3 Agnostic PAC-Learning

Remark 51. (Motivation for Agnostic PAC) Assuming that there exists f^* may be too strong, so we relaxed the notation, so we use the assumption that the joint distribution \mathcal{D} be a distribution over $\mathcal{X} \times \mathcal{Y}$ as now we are going to use:

$$L_{\mathcal{D}}(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}}[h(x) \neq y] := \mathcal{D}(\{(x,y) : h(x) \neq y\})$$

We will redefine the approximately correct notion.

Definition 7.12. (General Agnostic PAC) A hypothesis class \mathcal{H} is agnostic PAC learnable if there exists a function $m_{\mathcal{H}} : (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm A with the following properties: for every $\delta, \varepsilon \in (0,1)$ and $m > m_{\mathcal{H}}(\varepsilon, \delta)$:

$$\mathcal{D}^m \left(\left\{ S \in (\mathcal{X} \times \mathcal{Y})^m : L_{\mathcal{D}}(A(S)) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon \right\} \right) \geq 1 - \delta$$

Definition 7.13. (ε -Representation Sample) A training set S is called ε -representative if for all $h \in \mathcal{H}$ as we have:

$$|L_S(h) - L_{\mathcal{D}}(h)| \leq \varepsilon$$

Lemma 7.2. *Assume that a training set S is $\varepsilon/2$ -representative, then the average output of $\text{ERM}_{\mathcal{H}}(S)$ namely $h_S \in \arg \min_{h \in \mathcal{H}} L_S(h)$ satisfies:*

$$L_{\mathcal{D}}(h_S) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \varepsilon$$

Proof. For every $h \in \mathcal{H}$ as we have:

$$L_{\mathcal{D}}(h_S) \leq L_S(h_S) + \frac{\varepsilon}{2} \leq L_S(h) + \frac{\varepsilon}{2} \leq L_{\mathcal{D}}(h) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = L_{\mathcal{D}}(h) + \varepsilon$$

□

Definition 7.14. (Uniform Convergence) Let \mathcal{H} has the uniform convergence if there exists a function $m_{\mathcal{H}}^{\text{UC}} : (0,1)^2 \rightarrow \mathbb{N}$ such that for every $\varepsilon, \delta \in (0,1)$ and every distribution \mathcal{D} , and we have:

$$\mathcal{D}^m \left\{ \left(S \in Z^m : S \text{ is } \varepsilon\text{-representable} \right) \right\} \geq 1 - \delta$$

where Z is the domain and $m \geq m_{\mathcal{H}}^{\text{UC}} : (0,1)^2$

Corollary 7.1. *From the definition of uniform convergence, we can show that:*

- If \mathcal{H} has uniform convergence property with a function $m_{\mathcal{H}}^{\text{UC}}$ then \mathcal{H} is agnostic PAC learnable with sample complexity of

$$m_{\mathcal{H}}(\varepsilon, \delta) \leq m_{\mathcal{H}}^{\text{UC}}(\varepsilon/2, \delta)$$

This follows from the lemma above.

- We can show that $\text{ERM}_{\mathcal{H}}$ is successful against PAC learner for \mathcal{H} .

Theorem 7.8. Assume \mathcal{H} is finite, then \mathcal{H} is agnostic PAC learnable using $\text{ERM}_{\mathcal{H}}$ algorithm with:

$$m_{\mathcal{H}}(\varepsilon, \delta) = \left\lceil \frac{2 \log(2 |\mathcal{H}| / \delta)}{\varepsilon^2} \right\rceil$$

Comparing the reliable case generalization, the error will decrease in \sqrt{m} values as oppose to linear.

Proof. It suffices to show that \mathcal{H} has the uniform convergence property with:

$$m_{\mathcal{H}}^{\text{UC}}(\varepsilon, \delta) \leq \left\lceil \frac{2 \log(2 |\mathcal{H}| / \delta)}{\varepsilon^2} \right\rceil$$

To show that the uniform convergence, we need to show that:

$$\mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \varepsilon\}) \leq \delta$$

Using the union bound, we can see that:

$$\begin{aligned} \mathcal{D}^m(\{S : \exists h \in \mathcal{H}, |L_S(h) - L_{\mathcal{D}}(h)| > \varepsilon\}) &\leq \delta = \mathcal{D}^m\left(\bigcup_{h \in \mathcal{H}} \{S : |L_S(h) - L_{\mathcal{D}}(h)| \geq \varepsilon\}\right) \leq \delta \\ &\leq \sum_{h \in \mathcal{H}} \mathcal{D}^m(\{S : |L_S(h) - L_{\mathcal{D}}(h)| \geq \varepsilon\}) \leq \delta \\ &\leq 2 |\mathcal{H}| \exp(-2m\varepsilon^2) \end{aligned}$$

The last inequality is shown by Hoeffding inequality, setting the correct m , to finish the proof. \square

Remark 52. (Error Decomposition) Let $h_S = \text{ERM}_{\mathcal{H}}(S)$, we can decompose the risk as:

$$L_{\mathcal{D}}(h_S) = \mathcal{E}_{\text{app}} + \mathcal{E}_{\text{est}}$$

We have the following error:

- *Approximation Error:* $\mathcal{E}_{\text{app}} = \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$. How much risk do we need to restrict \mathcal{H} ? This doesn't depend on S , while it decreases with the complexity of \mathcal{H} increases.
- *Estimation Error:* $\mathcal{E}_{\text{est}} = L_{\mathcal{D}}(h_S) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$. It is the result of L_S being estimator of $L_{\mathcal{D}}$, while it decreases with size S but increase with complexity of \mathcal{H} .

This is bias and complexity: choosing $\mathcal{H}' \supset \mathcal{H}$ leads to decreases in \mathcal{E}_{app} while \mathcal{E}_{est} increases.