Stas Minsker

Intro    machine learning is the study of algorithms that can learn from data, gradually improving
their performance. large datasets

mathematical statistics (ST) is the science of making decision in the face of uncertainty.

small datasets

[Ex]    $Y = \alpha + \beta_1 [interesting] + \beta_2 [genre] + \beta_3 [Budget] + \cdots + \varepsilon$

↓        average rating (or?)

rating of a movie
    $\in [0, 100]$

Netflix: wants to predict rating of a movie  - ML

Disney:  wants to make a movie with high rating - ST
        needs to understand whether the model is statistically significant
        (hypothesis testing, etc. ···)

[Ex] Handwritten digits recognition  — ml

ml {
    unsupervised learning    "raw data", no label     clustering, learning representation
    supervised learning       labeled                  classification, prediction
    reinforcement learning    "multi-armed bandit problem"  exploration  v.s.  exploitation
}

Realizable case:

Binary classification   (X, Y) = (instance, label)
                                (observation, label)

[Ex] $x$ - image, $Y \in \{+1, -1\}$  (eg: "cat", "dog")
    $x \in S$ - set of all possible instances

Statistical learning:   we will assume that (X, Y) is random, in other words, it has a probability distribution P
        so we use language of probability theory

Supervised Learning:
    $(X, Y) \in S \times \{+1, -1\}$
    P is the distribution of (X, Y)
    i.e.  $P(A) =$ Probability $((X, Y) \in A)$
    $\Pi$ is the distribution of X
    Imposing the probabilistic model on (X, Y) takes as into realm of Statistical Learning
    Theory
    Goal: predict label Y based on the observation X
    The prediction rule is a function  $g: S \to \{-1, +1\}$

The quality of a _prediction rule_ g is measured by the classification / generalization error

$L(g) = \text{Prob}(Y \neq g(X))$   (one prediction)

The _training data_ is a sequence $(X_1, Y_1), (X_2, Y_2), \cdots, (X_n, Y_n)$ of i.i.d. pairs with distribution $P$.

a prediction rule

An algorithm takes training data as an input and outputs $\hat{g}_n = \hat{g}_n((X_1, Y_1), \cdots, (X_n, Y_n))$

In general, we will consider 2 scenarios:

1) "Realizable" learning : there exists $g \in G$ s.t. $Y = g^*(x)$ with probability 1.

2) "Agnostic" learning : there is no $g \in G$ s.t. $Y = g^*(x)$ with probability 1.

Realizable scenario :

    assume that the set $G$ of all possible classification rules is _finite_.

    By assumption, $\exists g^* \in G : Y = g^*(X)$ with prob 1

The _Empirical Risk_ Minimization principle :

  (training data) pick any $\hat{g}_n$ that agrees with the training data ( $\hat{g}_n(X_i) = Y_i$, $i = 1, \cdots, n$)

    Question : what is $L(\hat{g}_n)$ ?

          what is $\text{prob}(L(\hat{g}_n) > \varepsilon) \leq ?$ , given $\varepsilon > 0$

    Here, $L(g) = \text{Prob}(Y \neq g(X))$

    let $G_B =$ "bad" classification rules $= \{g \in G : L(g) > \varepsilon\}$

    $\text{prob}(L(\hat{g}_n) > \varepsilon) = \text{prob}(\hat{g}_n \in G_B)$

    Takes $g \in G_B$, if $\hat{g}_n = g$, $g(X_i) = Y_i$, $i = 1, 2, \cdots, n$

    $\text{prob}(g(X_1) \neq Y_1) > \varepsilon$

    $\text{prob}(g(X_1) = Y_1) \leq 1 - \varepsilon$

    $\Rightarrow \text{prob}(g(X_i) = Y_i, i = 1, \cdots, n) = P(g(X_1) = Y_1) \cdot P(g(X_2) = Y_2) \cdots P(g(X_n) = Y_n)$

                        $= \prod_{i=1}^{n} Pr(g(X_i) = Y_i)$

                        $\leq (1-\varepsilon)^n$

    We show that $\forall g \in G_B$, $Pr(\hat{g}_n = g) \leq (1-\varepsilon)^n \leq e^{-\varepsilon n}$

                    (since $1 - \varepsilon \leq e^{-\varepsilon}$)

    Remainder : Union Bound : $P(A \cup B) \leq P(A) + P(B)$

    $\therefore$ If $G_B = \{g_1, \cdots, g_k\}$, then $P(\hat{g}_n \in G_B) = P(\hat{g}_n = g_1 \text{ or } \hat{g}_n = g_2, \text{ or } \cdots, \text{ or } \hat{g}_n = g_k)$

                        $\leq P(\hat{g}_n = g_1) + P(\hat{g}_n = g_2) + \cdots + P(\hat{g}_n = g_k)$

                        $\leq k e^{-\varepsilon n}$

                        $\leq |G| e^{-\varepsilon n}$

If one requires that $P(L(\hat{g}_n) \leq \varepsilon) \geq 1 - \delta$

then $|G| e^{-\varepsilon n} \leq \delta$

$\Leftrightarrow n \geq \frac{\log \frac{|G|}{\delta}}{\varepsilon}$

E.g. if $G = 10^{24} = 2^{11}$, $\delta = 2^{-6}$, $\varepsilon = 0.01$, what is $n$? (最少有多少样本)

$n \geq \frac{16 \log 2}{0.01} = 160 \log 2$

Record some useful concepts:

ERM: pick a classifier that makes the smallest number of mistakes on observed data.

define: $L_n(g) = \frac{1}{n} \cdot \#\{1 \leq j \leq n : g(x_j) \neq Y_j\}$

Remark: By Law of Large Number, $L_n(g) \xrightarrow{P} L(g)$ since $\mathbb{E} L_n(g) = L(g)$

The ERM states that one pick $\hat{g}_n$ that minimize $L_n(g)$ over $g \in G$

PAC: ("Probably Approximately Correct") Learnability: Leslie Valiant '84

A class $H$ of hypothesis / binary classifiers is PAC learnable:

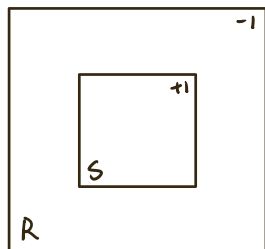if $\forall \varepsilon, \delta \in (0,1)$, any labeling function $g^*$, any distribution $\Pi$ of $x$,

$\exists$ algorithm $A$. (e.g. ERM) and a function $\underset{\text{called sample complexity}}{n = u(\varepsilon, \delta, G)}$ st. $\Pr(L(\hat{g}_n) > \varepsilon) \leq \delta$

推论: Any $\underset{\text{有}|G|}{\text{finite}}$ set of binary classifiers is PAC-learnable

if we take $n \geq \frac{\log(\frac{|G|}{\delta})}{\varepsilon}$, we can satisfy $\Pr(L(\hat{g}_n) > \varepsilon) \leq \delta$

Infinite set of classifiers

[Example A]

Area$(R) = 2$ (大方块)

Area$(S) = 1$

$g^*(x) = \begin{cases} +1 & x \in S \\ -1 & x \in R \backslash S \end{cases}$

let $G = \{\text{all binary function } g : R \to \{+1, -1\}\}$

Training Data: $(x_1, Y_1), \cdots, (x_n, Y_n)$

consider $\hat{g}_n(x) = \begin{cases} Y_i, & x = x_i \text{ for } i = 1, \cdots, n \\ -1, & \text{else} \end{cases}$

In particular, $\hat{g}_n$ is consistent with ERM (sample 里见过的都对)

But $L(\hat{g}_n) = \overset{\text{loss}}{\Pr}(Y \neq \hat{g}_n(x)) = \frac{1}{2} \to$ overfitting

(assume that $x$ is chosen uniformly from $R$) 选到点的概率是0 $\Rightarrow$ 几乎全是-1

$x \sim U(0,1)$ $\underset{\overset{h}{x} \overset{h}{}}{\overset{\text{samples}}{\mid\mid\mid\mid\mid}}$

$\Pr(X = x) \leq \Pr(X \in [x-h, x+h]) \ \forall h$

since $h$ can be as small as it wants to

$\therefore \Pr(X = x) = 0$

not all infinite |G| is not PAC learnable ?
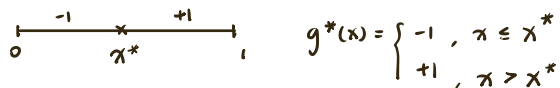
<span style="color:blue">overfitting</span> : $L(\hat{g}_n) \neq 0$ ?

if $G$ is too "large", then any <mark>algorithm</mark> (in particular ERM) will produce a classifier with large misclassification error.

If $G$ is too "large" → it is not PAC learnable

An <u>algorithm</u> A is a map that takes $G$ and $(X_i, Y_i)$, $i = 1, \cdots, n$ as input and outputs $\hat{g} \in G$.
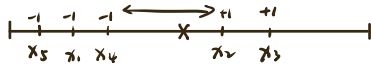
[Example B]  $x \in [0,1]$



$$g^*(x) = \begin{cases} -1 & , x \leq x^* \\ +1 & , x > x^* \end{cases}$$

$G = \{ g_y , y \in [0,1] \}$ ← 有定义域，不是 all classifier, 但仍 infinite classifiers

$$g_y(x) = \begin{cases} -1 & , x \leq y \\ +1 & , x > y \end{cases} \Big\} \text{只能分两段}$$

$(X_1, Y_1), \cdots, (X_n, Y_n)$ iid ― training data



claim : $G$ is PAC learnable ?

→ To show this, we need to estimate its sample complexity
  given $\varepsilon, \delta > 0$, if $n \geq \underline{n(\varepsilon, \delta)}$, $\exists A$ such that, given a sample of size $n$,
  A outputs $\hat{g}$ such that $L(\hat{g}) < \varepsilon$ with prob $\geq 1 - \delta$
  let's use ERM: specifically, let $\hat{x} = \max \{ x_j : Y_j = -1 \}$
  let $\hat{g}_n(x) = \begin{cases} -1 & , x \leq \hat{x} \\ +1 & , x > \hat{x} \end{cases}$

  $Pr(L(\hat{g}) > \varepsilon) \leq (1 - \varepsilon)^n$   (show this!)

[Example 1]  Task: identify counterfeit banknotes

We know that real banknotes  (a) color change under the light $\in [0, 1]$, with increment $0.1$

(b) red/blue fibers $\in [0, 100]$, with increment 1

realizable ← Assume that for every "real" banknote, (a)(b) belong to a specific range, vice versa

Assume that you have 100 banknotes, known to be real or not
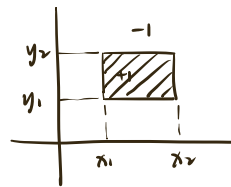
→  $X$: banknote,  $Y$: label { real / fake }

$x \in S$   Domain set $= \{ (x, y) : x \in [0 : 0.1 : 1], y \in [0 : 1 : 100] \}$

Training data: 100 banknotes

Hypothesis class $G$:  $g : S \rightarrow \{ +1, -1 \}$

$g(x) = \begin{cases} +1, & x \in [x_1, x_2], y \in [y_1, y_2] \\ -1, & else \end{cases}$

$|G| \leq \binom{11 \times 100}{2} \div 2$  — still manageable



Assume that we want a classifier that makes at most 5% mistake.

What is the probability that you will get such a classifier from a sample size 100.

We proved that

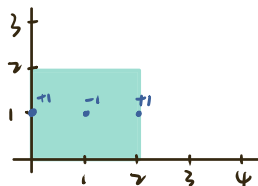$$Pr( L(\hat{g}_{ERM}) > 0.05) \leq \frac{10^8}{2} (1 - 0.05)^{100}$$  ↗ 代表 $|G|$

$$= \frac{10^8}{2} \cdot \frac{95^{100}}{10^{200}}$$   prob 小于1万里4几 东西?

[Example 2]  label $(x_1, x_2)$,  $x_1, x_2 \in \mathbb{I}$,  $0 \leq x_1 \leq 4$,  $0 \leq x_2 \leq 3$

$G = \{$ rectangle with vertices $(x_1, x_2) \in [0,4] \times [0,3]) \}$

Training set:

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 1 | 1 |
| 1 | 1 | -1 |
| 2 | 1 | 1 |



let $g : [0, 2] \times [0, 2]$

$L_n(g) = \frac{\#\{1 \leq i \leq n, Y \neq g(x_i)\}}{n} = \frac{1}{3}$  ← empirical risk

$(n=3)$

ERM: line segment $(0, 1)$ to $(2, 1)$

→ Agnostic learning ( no perfect classifier)

An algorithm $A$ is a mapping from training data $(x_i, y_i)_{i=1}^{n}$ to the class $G$ of binary classifiers

Theorem: Assume that $S$ is finite. Let $(x_1, y_1), \cdots, (x_n, y_n)$ be the training data such that $n \leq \frac{|S|}{2}$.

Then for any algorithm $A$, $\exists$ some distribution $\pi$ over $S$ and $g^*: Y = g^*(x)$ but $Pr(\hat{g}(x) \neq g^*(x)) \geq \frac{1}{8}$ with prob $\frac{1}{8}$ where $\hat{g} = A((x_i, y_i)_{i=1}^{n})$   P(习得算法犯错概率 $> \frac{1}{8}$) $> \frac{1}{8}$

Proof:   let $\mathcal{X} = \{(x_1, y_1), \cdots, (x_n, y_n)\}$

consider $\max_{g^*, \pi} E_{\mathcal{X}} E_{x \sim \pi} I\{\hat{g}(x) \neq g^*(x)\}$

$\underbrace{Pr(\hat{g}(x) \neq g_*(x))}_{} \geq \underset{\text{Want}}{c > 0}$   $\frac{1}{8}$

pick $g^{*(1)}, g^{*(2)}, \cdots, g^{*(k)}_*$

$\max(\quad) \geq \frac{1}{k} \sum_{j=1}^{k} E_{\mathcal{X}} E_{x \sim \pi} I\{\hat{g}(x) \neq g^{*(j)}_*(x)\}$   最大值 $\geq$ expected

consider a random $g^*$ s.t. $\forall x \in S$, $g_*(x) = \begin{cases} +1 & \text{with prob } \frac{1}{2} \quad \text{independently.} \\ -1 & \text{with prob } \frac{1}{2} \end{cases}$

call this distribution $Q$ $\nearrow$

$\max_{g^*, \pi} E_{\mathcal{X}} E_{x} I\{\hat{g}(x) \neq g_*(x)\} \geq E_{g^* \sim Q} E_{\mathcal{X}} E_{x} I\{\hat{g}(x) \neq g^*(x)\}$

$= E_{\mathcal{X}} E_{x} E_{g^* \sim Q} I\{\hat{g}(x) \neq g^*(x)\}$

$\geq E_{\mathcal{X}} E_{x} \frac{1}{2} I\{x \notin \{x_1, \cdots, x_n\}\}$

$E_{g^* \sim Q} I\{\hat{g}(x) \neq g_*(x)\} = \begin{cases} 0 & \text{if } x \in \{x_1, \cdots, x_n\} \quad \text{如果见过,那么会是对的} \\ \frac{1}{2} & \text{if } x \notin \{x_1, \cdots, x_n\} \quad \text{如果没见过,那么有一半概率会对} \end{cases}$

③ $\begin{cases} \geq E_{\mathcal{X}} \frac{1}{2} \cdot \frac{1}{2} \\ = \frac{1}{4} \end{cases}$

[lemma] let $Z$ be a r.v. such that $0 \leq Z \leq 1$, then $Pr(Z \geq \delta) \geq E(Z) - \delta$

proof: $EZ = EZ \cdot I\{Z \leq \delta\} + EZ \cdot I\{Z > \delta\}$

$\leq \delta + Pr(Z > \delta)$

$\therefore$ Applying the lemma, we get that

$\max_{g^*} Pr_{\mathcal{X}} (Pr(\hat{g}(x) \neq g_*(x)) \geq \frac{1}{8}) \geq \frac{1}{4} - \frac{1}{8} = \frac{1}{8}$

$(\because E_{\mathcal{X}}(Pr(\hat{g}(x) \neq g^*(x))) = \frac{1}{4}$

[Lemma] let $Z$ be a r.v. s.t. $\text{Var}(Z) < \infty$

Then, $EZ = \underset{z \in \mathbb{R}}{\text{argmin}}\ E(Z-z)^2$

Proof: let $f(z) = E(Z-z)^2$

Then, $f'(z) = 2E(Z-z) = 0 \iff z = EZ$

Finally, $f''(z) = 2 \Rightarrow EZ$ is the minimizer

(弄口同上)

Now, let $Z, W$ be such that $\text{Var}(Z)$ and $\text{Var}(W)$ are finite.

Then $E[Z | W=y] = \underset{z \in \mathbb{R}}{\text{argmin}}\ E_{Z|W=y} (Z-z)^2$

Clearly, $z = z(y)$ above. Therefore, $E[Z | W=y]$ is a function of $W$ that minimizes $E[Z - f(W)^2]$ over all functions $f$.

[Exercise] let $y(w) = E[Z|W]$. Prove that for any function $g$, $E(Z - y(w)) \cdot g(w) = 0$

let $h(Z)$ be an arbitrary function of $Z$ s.t. $Eh^2(Z) = \infty$

then $E[(w - \hat{f}(Z) \cdot h(Z)] = 0$ where $\hat{f}(Z) = E[w|Z]$

$h(z)$ $\hat{f}(Z)$

$f(z)$

$\langle x, y \rangle = E(x, y)$

## Bayes Classifier

$$y(x) = E[Y \mid X = x] \qquad Y \in \{+1, -1\}$$

**Theorem**  let $S$ be a finite set. $X$ has (discrete) distribution $\Pi$ over $S$. Then the best possible binary classifier is given by $g^*(x) = \text{sign}(E(Y \mid X = x))$

$g^*(x)$ is known as **Bayes classifier**

**Proof:**  let $g: S \to \{\pm 1\}$ be arbitrary

Then $Pr(Y \neq g(x)) = \sum_{x \in S} Pr(Y \neq g(x) \mid X = x) \cdot \Pi(x)$    *probability of $X = x$*

$Pr(Y=1 \mid X=x) +$
$Pr(Y=-1 \mid X=x) = 1$
$\begin{cases} Pr(Y=1 \mid X=x) = \frac{1+y(x)}{2} \\ Pr(Y=-1 \mid X=x) = \frac{1-y(x)}{2} \end{cases}$ $\iff$ $Pr(Y=t \mid X=x) = \frac{1+ty(x)}{2}$, $t \in \{\pm 1\}$
$\xrightarrow{\text{对正确率}}$ $Pr(Y \neq t \mid X=x) = \frac{1-ty(x)}{2}$,

Then $\underline{Pr(Y \neq g(x))} = \sum_{x \in S} \frac{1-g(x)y(x)}{2} \cdot \Pi(x)$

*想要 minimize*

$$\geq \sum_{x \in S} \frac{1 - |y(x)|}{2} \Pi(x)$$

Equality is achieved when $g(x)y(x) = |y(x)|$ for all $x \in S$

$$\iff \underline{g(x) = \text{sign}(y(x))} \to \text{Bayes Classifier}$$

## Bayes Risk

$$L^* = L(g^*) = \sum_{x \in S} \left( \frac{1 - |y(x)|}{2} \right) \Pi(x)$$

**[Example]**  5 cards are drawn at random. 2 are reviewed

$$Y = \begin{cases} 1, & 5 \text{ cards contain an Ace} \\ -1, & \text{otherwise} \end{cases}$$

Find the Bayes classifier and its risk.

**Remark:**  The risk of the Bayes classifier is called the Bayes risk: $L^* = Pr(Y \neq g^*(x))$

**Solution:**  $S \in \{1, 0\} \to \begin{cases} x=1 \text{ if the pair of cards have at least 1 Ace} \\ x=0 \text{ otherwise} \end{cases}$

$y(x) = E[Y \mid X = x] = 1 \cdot Pr(Y=1 \mid X=x) + (-1) \cdot Pr(Y=-1 \mid X=x)$

$Pr(Y=-1 \mid x=1) = 0$

$Pr(Y=-1 \mid x=0) = \frac{\binom{46}{3}}{\binom{50}{3}} \approx 0.77$

$Pr(Y=1 \mid x=1) = 1$
$Pr(Y=1 \mid x=0) = 1 - \frac{\binom{46}{3}}{\binom{50}{3}}$

$\underset{x=0}{y(0)} = 1 \cdot (1-0.77) + (-1) \cdot 0.77 = -0.54$
$\underset{x=1}{y(1)} = 1 \cdot (1-0) + (-1) \cdot 0 = 1$

$g^*(x) = \text{sign}(y(x)) = \begin{cases} 1, & x=1 \\ -1, & x=0 \end{cases}$

$$P(x=1) = 1 - P(\text{both cards are not Aces}) = 1 - \frac{48}{52} \cdot \frac{47}{51} = 0.15$$

$$P(x=0) = 1 - 0.15 = 0.85$$

$$L* = \frac{1-|y(0)|}{2} \cdot \Pi(0) + \frac{1-|y(1)|}{2} \cdot \Pi(1)$$

$$= \frac{1-0.54}{2} \cdot 0.85 + \frac{1-1}{2} \cdot 0.15$$

$$\approx 0.2$$

## Agnostic PAC - learnability

(APAC(L))

A class $G$ of binary classifiers is agnostic PAC-learnable if $\exists\, m(G; \varepsilon, \delta)$ and an algorithm $A$ such that $\forall\, \varepsilon, \delta > 0$, any distribution $P$ of $(X, Y)$, $L(\hat{h}_n) \leq \min\limits_{g \in G} L(g) + \varepsilon$ with prob $\geq 1 - \delta$ as long as $n \geq m(G; \varepsilon; \delta)$.

<span style="color:blue">↑ output of A    $L(g^*)$</span>

Next goal: understand which classes are PAC learnable. We will start with finite classes, and then will study infinite classes.

key idea: concept of <u>uniform closeness</u> of the true and empirical risks.

Definition   The training set $X = (X_1, Y_1), \ldots, (X_n, Y_n)$ is called <u>$\varepsilon$-representative</u> if   <span style="color:blue">$L(g) = Pr(Y \neq g(x))$</span>

$$\forall g \in G, \ |L_n(g) - L(g)| \leq \varepsilon$$

<span style="color:blue">$L_n(g) = \frac{1}{n} \cdot \#\{1 \leq j \leq n ; g(x_j) \neq y_j\}$</span>

|empirical risk - true risk| $\leq \varepsilon$    <span style="color:blue">$= \frac{1}{n} \sum\limits_{j=1}^{n} I\{Y_j \neq g(x_j)\}$</span>

[lemma]   Assume that $X$ is $\frac{\varepsilon}{2}$ representative, and let $\hat{g}_n$ be the minimizer of the empirical risk:

$$\hat{g}_n = \arg\min_{g \in G} L_n(g). \text{ Then}$$

$$L(\hat{g}_n) = Pr(Y \neq \hat{g}_n(x) \mid x) \leq \min_{g \in G} L(g) + \varepsilon$$

[proof]   let $\bar{g} = \arg\min_{g \in G} L(g)$. Then

$$L(\hat{g}_n) = L_n(\hat{g}_n) + L(\hat{g}_n) - L_n(\hat{g}_n)$$

$$\leq \underbrace{L_n(\bar{g})}_{\color{blue}{L_n(\hat{g}_n) \leq L_n(\bar{g})}} + \max_{g \in G} |L(g) - L_n(g)|$$

<span style="color:blue">已是empirical<br>最优解</span> $\leq L(\bar{g}) + \frac{\varepsilon}{2} + |L_n(\bar{g}) - L(\bar{g})|$

$$\leq L(\bar{g}) + \frac{\varepsilon}{2} + \frac{\varepsilon}{2}$$

$$\leq L(\bar{g}) + \varepsilon$$

## Bias-Complexity Tradeoff

This concept refers to the following error decomposition:

let $\hat{g}$ be the output of a learning algorithm $A$ given the training data $\mathcal{X} = (X_1, Y_1), \cdots, (X_n, Y_n)$

Then
$$L(\hat{g}) = Pr(Y \neq \hat{g}(X) | \mathcal{X})$$
$$= L(\hat{g}) - \min_{g \in G} L(g) + \underbrace{\min_{g \in G} L(g)}_{\geq \text{ Bayes Risk } L^*}$$

If $G$ is large, $\min\limits_{g \in G} L(g)$ is small. But $L(\hat{g}) - \min\limits_{g \in G} L(g)$ is large

## Finite Classes are agnostic PAC learnable

Question: What is the smallest sample size sufficient to guarantee that it is $\varepsilon$-representative with probability at least $1 - \delta$?

Assume $|G| < \infty$, then $Pr(\forall g \in G, |L_n(g) - L(g)| \leq \varepsilon)$
$$= 1 - Pr(\exists g \in G, |L_n(g) - L(g)| > \varepsilon)$$

*why can apply union bound ②*

$$Pr(\exists g \in G, |L_n(g) - L(g)| > \varepsilon)$$
$$= Pr\left( \bigcup_{g \in G} \{ |L_n(g) - L(g)| > \varepsilon \} \right)$$

*we need to show the measure of $L_n(g)$ is concentrated around its expected value*

$$\leq \sum_{g \in G} Pr(|L_n(g) - L(g)| > \varepsilon)$$
$$\leq |G| \max_{g \in G} Pr(|L_n(g) - L(g)| > \varepsilon)$$

apply chebyshev's inequality:  $P(|x - E(x)| \geq k) \leq \dfrac{\sigma^2}{k^2}$

Fix $g \in G$, $Pr\left( \left| \frac{1}{n} \sum_{j=1}^{n} I\{ Y_j \neq g(X_j) - L(g) \} \right| > \varepsilon \right)$

$\qquad L(g) = E L_n(g)$

$\qquad\quad Z_j = I\{ Y_j \neq g(X_j) \}$

$\qquad\quad Z_1, \cdots, Z_n$ are i.i.d. $\to \mathbb{Z}$

$Pr\left( \left| \frac{1}{n} \sum Z_j - E\mathbb{Z} \right| > \varepsilon \right) \leq \dfrac{Var(\frac{1}{n} \sum Z_j)}{\varepsilon^2}$

$Var\left( \frac{1}{n} \sum Z_j \right) = \sum Var\left( \frac{1}{n} Z_j \right) = n \cdot \frac{1}{n^2} Var(\mathbb{Z}) = \dfrac{Var(\mathbb{Z})}{n}$

$\leq |G| \dfrac{Var(\mathbb{Z})}{n \varepsilon^2} \to P(\mathcal{X} \text{ is not } \varepsilon\text{-representative})$

$$\dfrac{|G| Var(\mathbb{Z})}{n \varepsilon^2} \leq \delta \quad \Rightarrow \quad n \geq \dfrac{|G| Var(\mathbb{Z})}{\delta \varepsilon^2} \qquad \overset{?}{\circ} \text{ Bad Bound}$$

$Pr(Z=1) = Pr(Y \neq g(x)) = Pg$

$Pr(Z=0) = 1 - Pg$

$Var(Z) = Pg(1-Pg) \leq \frac{1}{4}$ → suppose Bernoulli distribution

$f(x) = x(1-x)$



$\therefore Pr(* \text{ is not } \varepsilon - \text{representative }) \leq |G| \cdot \frac{1}{4n\varepsilon^2}$

[Exercise] $|G| = 1000$

we want $*$ $\varepsilon$-representative with prob at least $0.9$ and $\varepsilon = 0.1$ $(\delta = 0.1)$

$|G| \frac{1}{4n\varepsilon^2} \leq 0.1 = \delta$

$n \geq \frac{|G|}{4\delta\varepsilon^2} = \frac{10^3}{4 \times 10^{-3}} = 250,000$ 很大, bad estimation

$\therefore *$ is $\varepsilon$-representative with prob $1-\delta$

$\Leftrightarrow \hat{g}_n$ obtained by ERM satisfy $L(\hat{g}_n) \leq \min_{g \in G} L(g) + 2\varepsilon$ with prob $1-\delta$

## Hoeffding's Inequality

**Lemma 4.5** (Hoeffding's Inequality). *Let $\theta_1, \ldots, \theta_m$ be a sequence of i.i.d. random variables and assume that for all $i$, $\mathbb{E}[\theta_i] = \mu$ and $\mathbb{P}[a \leq \theta_i \leq b] = 1$. Then, for any $\epsilon > 0$*

$$\mathbb{P}\left[\left|\frac{1}{m}\sum_{i=1}^{m}\theta_i - \mu\right| > \epsilon\right] \leq 2\exp\left(-2m\epsilon^2/(b-a)^2\right).$$

apply Hoeffding's inequality to question:

let $\theta_i$ be the random variable $L_n(g)$ , $\theta_1, \cdots, \theta_n$ are i.i.d.

(∵ $L_n(g)$ 选取的 sample data 也是 i.i.d. 的)

$L_n(g) = \frac{1}{n}\sum_{i=1}^{n}\theta_i$ , $L(g) = EL_n(g) = \mu$

$\therefore P(|\frac{1}{n}\sum_{i=1}^{n}\theta_i - \mu| > \varepsilon) \leq 2e^{-2n\varepsilon^2}$

$\therefore Pr(\exists g \in G, |L_n(g) - L(g)| > \varepsilon) \leq \sum_{g \in G} 2e^{-2n\varepsilon^2} = 2|G|e^{-2n\varepsilon^2}$

$\therefore n \geq \frac{\log(\frac{2|G|}{\delta})}{2\varepsilon^2}$ then $Pr(\exists g \in G, |L_n(g) - L(g)| > \varepsilon)$

对于上面 E.x. , $n \geq \frac{1}{2} \cdot 100 \cdot \log(\frac{2 \cdot 10^3}{0.1}) \leq 2000$

more reasonable estimation

[Exercise] $S = \{0, 1, 2, 3, 4\}$  $X$ is binomial $B(4, \frac{1}{2})$, $Y \in \{+1, -1\}$

$P(Y=1 | X=x) = \frac{1}{2}$  (label is random guess)

$x \sim B(n, p)$
$Pr(x=k) = \binom{n}{k} p^k (1-p)^{n-k}$

consider $g(x) = \begin{cases} 1, & x \text{ is even} \\ -1, & x \text{ is odd} \end{cases}$

interpretation: $L(g) = Pr(Y \neq g(x)) = \frac{1}{2}$

$y(x) = E[Y | X = x] = 0$

$g^*(x) = \text{sign}(y(x)) = \text{either } +1 \text{ or } -1$

Now, $P(Y=1 | X=x) = \begin{cases} \frac{3}{4} & x = 0, 1, 2 \\ \frac{1}{4} & x = 3, 4 \end{cases}$

$\Rightarrow P(Y=-1 | X=x) = \begin{cases} \frac{1}{4} & x = 0, 1, 2 \\ \frac{3}{4} & x = 3, 4 \end{cases}$

$L(g) = Pr(Y \neq g(x)) = \sum Pr(Y \neq g(x) | X=x) \cdot Pr(X=x)$

$= \underbrace{(\frac{1}{2})^4 \times \frac{1}{4}}_{x=0} + \underbrace{\binom{4}{1} \times (\frac{1}{2})^4 \times \frac{3}{4}}_{x=1} + \underbrace{\binom{4}{2} \times (\frac{1}{2})^4 \times \frac{1}{4}}_{x=2} + \underbrace{\binom{4}{3} \times (\frac{1}{2})^4 \times \frac{1}{4}}_{x=3} + \underbrace{\binom{4}{4} \times (\frac{1}{2})^4 \times \frac{3}{4}}_{x=4}$

[Exercise]  $S = \mathbb{R}^2$, $G = \{g_r, r > 0\}$

$g_r(x) = \begin{cases} +1, & \|x\| \leq r \\ -1, & \|x\| > r \end{cases}$   Assume realizability
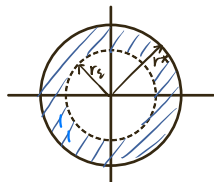
Show that $G$ is PAC-learnable


$Y = g_{r*}(x)$

→ show 2 things: (a) An algorithm $A \to$ ERM

(b) $m(\varepsilon, \delta)$ s.t. $\hat{g} = \hat{A}((x_1, Y_1), \cdots, (x_n, Y_n))$ satisfies $Pr(L(\hat{g}) \geq \varepsilon) \leq \delta$

ERM outputs $\hat{g}$ that minimizes $\#\{1 \leq j \leq n ; Y_j \neq g(x_j)\}$

$\therefore L_n(\hat{g}) = 0$

$\hat{r} = \min\{r > 0 : \|x_j\| < r \Leftrightarrow Y_j = +1\}$



let $r_\varepsilon$ be s.t. $Pr(r_\varepsilon \leq \|x\| \leq r_*) = \varepsilon$

(i) $Pr(\|x\| \leq r_*) < \varepsilon \Rightarrow$ Circle $(r = r_*)$ 面积小于 $\varepsilon \Rightarrow$ 所有 classifier loss 都小于 $\varepsilon$

(ii) $Pr(\|x\| \leq r_*) > \varepsilon$

$P(L(\hat{g}_r) \geq \varepsilon) = Pr(\hat{r} < r_\varepsilon) = Pr(\text{there are no instances with label } +1 \text{ in the ring}$
between circle of radius $r_\varepsilon$ and $r_*)$

$= Pr(\|x_j\| > r_* \text{ or } \|x_j\| < r_\varepsilon \text{ for all } j)$

$= \prod_{j=1}^{n} Pr(\|x_j\| > r_* \text{ or } \|x_j\| < r_\varepsilon)$
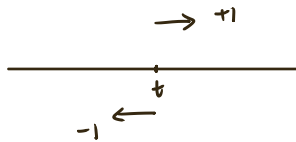
$$= (1-\zeta)^n \le e^{-\zeta n}$$

$$(\because \Pr(r_{\zeta} \le \|x\| \le r^*) = \zeta)$$

To have $e^{-\zeta n} \ge \delta$, $n \ge \dfrac{\log(\frac{1}{\delta})}{\zeta}$ $\qquad$ ($\Leftarrow$) finite class: $n \ge \dfrac{\log(\frac{|G|}{\delta})}{\zeta}$

[Exercise] let $S = R$, $G = \{g_t, t \in R\}$, $g_t(x) = \begin{cases} +1, & x \ge t \\ -1, & x < t \end{cases}$

Prove that $G$ is PAC-learnable assume realizability



Vapnik - Chervonenkis

Question: Which classes $G$ are agnostic PAC learnable?

observation: let $(X_1, Y_1), \cdots, (X_n, Y_n)$ be the training data, and $G$ is the concept class

consider $G_c = \{(g(x_1), \cdots, g(x_n)), g \in G\}$

note that $|G_c| \le 2^n$

from exercise, we have $g_r$ (圆形 classifier)

$g_r \quad \|x_1\|, \cdots, \|x_n\|$

$\quad i_1, \cdots, i_n$ is a permutation such that $\|x_{i_1}\| \le \|x_{i_2}\| \le \cdots \le \|x_{i_n}\|$

$\Rightarrow \{(g_r(x_{i_1}), \cdots, g_r(x_{i_n})), g_r \in G\}$

$\quad$ can tell $\quad Y_{ij} = \begin{cases} +1 \Rightarrow Y_{ik} = +1 \ \forall \ k \le j \\ -1 \Rightarrow Y_{ik} = -1 \ \forall \ k \ge j \end{cases}$

$\Rightarrow$ at most $(n+1)$ vectors

$\quad \rightarrow$ what makes it PAC-learnable

Let $(X_1, Y_1) \dots (X_n, Y_n)$ is the truining data
$G$ is the concept class
$C = (X_1, \dots X_n)$. The set $G_c = \{ (g(X_1) \dots g(X_n)), g \in G \}$
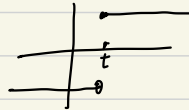
$G_c$ : restriction of $G$ onto $C$.
$|G_c| \leq 2^n$
  If $|G_c| = 2^n$, we will say $G$ __shatters__ $C$
__Remark__ : $C$ can be an arbitrary finite set


__Ex__ $G = \{ g_t, t \in \mathbb{R} \}$
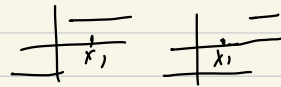  $g_t(x) = \begin{cases} +1, & x \geq t \\ -1, & x < t \end{cases}$



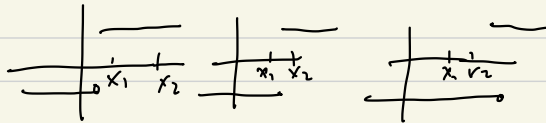Shift $t$ to get $\pm 1$ for $X_1$

$C = \{X_1\}$    $G_c = \{ g_t(X_1), t \in \mathbb{R} \}$
              $= \{ +1, -1 \}$      shatters



$C = \{X_1, X_2\}$



$(+1, +1)$     $(-1, +1)$     $(-1, -1)$     $(+1, -1)$
$G_c = \{ (1,1), (-1,-1), (-1, +1) \}$ no shatter


__Def__ $VC(G)$ — the Vapnik-Chervenkis dimension
  of $G$ — is the __largest__ $d$ such that $\exists \{X_1, \dots X_d\}$
  that is shattered by $G$

__Remark__ : $VC(G) = d \iff \begin{cases} \exists \{X_1, \dots, X_d\} \text{ shattered by } G \\ \text{any set } \{X_1, \dots X_{d+1}\} \text{ is not shattered by } G. \end{cases}$

<u>Remark:</u> we will prove the "Fundamental theorem of PAC learning"

    $G$ is agnostic PAC learnable $\iff VC(G) < \infty$

<u>Ex 2</u> $S$ is $\boxed{\text{infinite}}$

$$G = \{ g_T , T \subseteq S, |T| < \infty \}$$

$$g_T(x) = \begin{cases} 1, & x \in T \\ -1, & x \notin T \end{cases}$$

Then $VC(G) = \infty$

<u>Solution:</u> for any $d \geq 1$, we need to find $\{X_1, \ldots, X_d\} = C$
s.t. $|G_c| = 2^d$

Let's take any $\{x_1 \ldots x_d\}$, $G_c = \{ (g_T(x_1) \ldots g_T(x_d)), T \subseteq S, |T| < \infty \}$
$W = \{+1, -1\}^d$, $J = \{ j ; w_j = +1 \}$
$\underbrace{(+1}_{1}, -1, -1, -1, +1, \ldots \underbrace{+1}_{d} \}$
$\phantom{(+1, -1, -1,}\underset{4}{}$
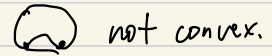
Take $T = \{ X_j, j \in J \}$ $\implies (g_T(x_1) \ldots g_T(x_d)) = W$

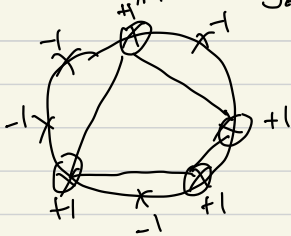<u>Ex 3</u> $G = \{ g_A, A \text{ is convex set} \}$

    $S = \mathbb{R}^2$ $\qquad g_A(x) = \begin{cases} 1, & x \in A \\ -1, & x \notin A \end{cases}$



convex

not convex.

$VC(G) = \infty$

<u>Ex 4</u> $G$ is $\boxed{\text{finite}}$ $\qquad VC(G) < \infty$

$(X_1, \ldots, X_d) = C$

$G_c = \{ g(X_1), \ldots g(X_d) \}, g \in G \}$

$|G_c| \leq |G|$

If $VC(G) = d \implies |G_c| = 2^d$ $\qquad |G| \geq 2^d \implies d \leq \lfloor \log_2 |G| \rfloor + 1$

                                            (→ integer!)

<u>lemma 1</u>: Let G be a concept class of infinite VC dimension, then G is <u>not</u> PAC-learnable

<u>Proof</u>: G is PAC-learnable if $\forall \varepsilon, \delta > 0$,
$\exists A$ - an algorithm and $m = m(\varepsilon, \delta)$  m # of training data must form.
then if $(X_1, Y_1)...(X_n, Y_n)$ is the training data and
$n \geq m(\varepsilon, \delta)$ then $Pr(L(\hat{g}) \geq \varepsilon) \leq \delta$, $\hat{g} = A((X_1, Y_1)...(X_n, Y_n))$
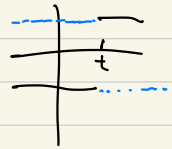
$VC(G) = \infty \Rightarrow \forall m \geq 1, \exists \{X_1, ..., X_n\} \subseteq S$, such that G shatters $\{X_1...X_n\}$
$\Rightarrow$ for any $\bar{z} \in \{+1, -1\}^m, \exists g \in G$ st $\bar{z} = (g(X_1), ..., g(X_m))$
$$G_c = \{+1, -1\}^m$$

Take $\varepsilon = \frac{1}{8}, \delta = \frac{1}{10}$
Take N arbitrarily large. Find $C = \{X_1, ..., X_N\}$ shattered by G
By no free lunch theorem for any A, $\max\limits_{g*} Pr(L(\hat{g}) \geq \frac{1}{8}) \geq \frac{1}{8}$
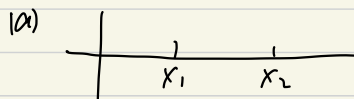
Proved by contradiction.

<u>EX 1</u>: $G = \{g_t^+, g_t^-, t \in \mathbb{R}\}$
$S = \mathbb{R}$, $g_t^+ = \begin{cases} 1 & x > t \\ -1 & x < t \end{cases}$   $g_t^- = \begin{cases} 1 & x < t \\ -1 & x \geq t \end{cases}$



Then $VC(G) = 2$   (Prove this)
(a) Find two points $\{X_1, X_2\}$ that are shattered.
(b) Show no set of 3 points are shattered.

(a)


$X_1$   $X_2$

$(+1, +1) \rightarrow g_t^+$ $t < X_1$
$(+1, -1) \rightarrow g_t^-$, $t \in (X_1, X_2)$
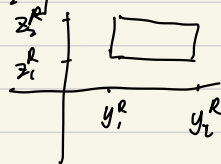$(-1, +1) \rightarrow g_t^+$, $t \in (X_1, X_2)$
$(-1, -1) \rightarrow g_t^-$, $t > X_2$

(b)  $(-1, +1, -1)$
$(+1, -1, +1)$   nope

Ex 2: $S = \mathbb{R}^2$, $G = \{g_R, R$ is axis-aligned Rectangle$\}$

$g_R = \begin{cases} +1, & X_1 \in [y_1^R, y_2^R], X_2 \in [z_1^R, z_2^R] \\ -1, & \text{else} \end{cases}$

(a)

$X_1$

$\circ K_4$   just need one example

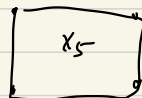$X_2$

$X_3$

(b) no 5 points are shattered

WLOG, assume that $X_1$ has the largest coordinate.
$X_2$ has the smallest $X$ coordinate
$X_3$ has the largest $y$ coordinate
$X_4$ has the smallest $y$ coordinate
$X_5$ will be inside the rectangle with $+$
which is impossible.

$X_5$

---

Theorem (R. Dudley) {not in textbook}
Let $L$ be a finite-dimensional space of function $f : S \rightarrow \mathbb{R}$. Consider
$C_f = \{\{x : f(x) > 0\}, f \in L\}$ and $\overline{C_f} = \{\{x : f(x) \geq 0\}, f \in L\}$.
Let $G = \{I_{C_f} - I_{C_f^c}, f \in F\}$   $\overline{G} = \{I_{\overline{C_f}} - I_{\overline{C_f}^c}, f \in F\}$
Then $VC(G) = VC(\overline{G}) = dim(L)$

Finite dim: Ex   $L = \{\langle a, x \rangle + b, a \in \mathbb{R}^d, b \in \mathbb{R}\}$
$a \in \mathbb{R}^d \Rightarrow a = a_1 e_1 + \dots + a_d e_d$
$e_j = (0, \dots, 0, 1, 0, \dots, 0)$
       $\uparrow$
       $j$
$\langle a, x \rangle + b = a_1 \langle e_1, x \rangle + \dots + a_d \langle e_d, x \rangle + b$
$dim(L) = d + 1$
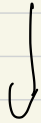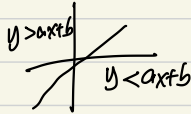         $\downarrow$
         constant

**Example**  Polynomials of degree at most $d$.
$$f(x) = a_1 x^d + a_2 x^{d-1} + \ldots + a_d x^2 + a_{d+1}$$   $\dim(L) = d+1$

$y = ax + b$        $\{ (x, y) : y - ax - b > 0 \}$
$y - ax - b = 0$

$y > ax + b$

$y < ax + b$

$a_1 x + a_2 y + b = 0$       $\dim = 3 = VC(G) = VC(\bar{G})$



rectangle formation:
**intersection of**
4 half subspaces
$\dim x 4 = 12$ (upper bound?)
$\frac{1}{3}$

**Proof** for (R. Dudley) :  $VC(G) \leq \dim(L)$
let $\dim(L) = d$, we need to show no set of $d+1$ points is shattered by $G$.
Take $\{x_1, \ldots, x_{d+1}\}$, Consider $T(f) = (f(x_1) \ldots f(x_{d+1}))$   $d+1$
$\qquad\qquad\qquad T(\alpha f + \beta g) = \alpha T(f) + \beta T(g)$
Note that $\dim(\text{Image}(T)) \leq d$ because $\dim(L) = d$ and linear maps don't increase dimension.
$\exists w \in \mathbb{R}^{d+1}$ such that $w \perp \text{Image}(T)$
and $w \neq 0$, Hence, if $w = (w_1, w_2, \ldots, w_{d+1})$



$\Rightarrow \exists j$ st $w_j < 0$ (if $w_j > 0$, take $-w$ instead)
$A_- = \{1 \leq j \leq n : w_j < 0\}$    $A_+ = \{1 \leq j \leq n, w_j \geq 0\}$
Assume that $\{x_1, \ldots, x_{d+1}\}$ is shattered by $G$
Since every $g \in G$ is sign $(f)$ for some $f \in F$.
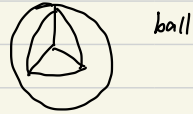$\exists f \in F$ st $f(x_j) > 0$, $j \in A_-$, $f(x_j) < 0$, $j \in A_+$

On one hand, since $w \perp \text{Im}(T)$
$$\sum_{j=1}^{d+1} w_j f(x_j) = 0$$

On the other hand, $\sum_{j=1}^{d+1} w_j f(x_j) = \underbrace{\sum_{j \in A_-} w_j f(x_j)}_{<0} + \underbrace{\sum_{j \in A_+} w_j f(x_j)}_{\leq 0}$  $< 0$

Contradiction $\Rightarrow \{x_1, \dots x_{d+1}\}$ cannot be shattered by $G$.

---

Example: $B_d(x, r) = \{y \in \mathbb{R}^d, \|y - x\|_2 \leq r\}$   ball
$G = \{g : \mathbb{R}^d \to \{+1, -1\}\}$, where
$g = g_{x,r}$, and $g_{x,r}(y) = \begin{cases} 1, & y \in B_d(x,r) \\ -1, & \text{else} \end{cases}$

Then $VC(G) \leq d+2$
Express definition of a ball as $f(y) \geq 0$ for $f \in L$, where $\dim(L) = d+2$

norm $= \sqrt{\sum_{j=1}^{d} (y_j - x_j)^2} \leq r \iff \sum_{j=1}^{d} (y_j - x_j)^2 \leq r^2 \iff \sum_{j=1}^{d} (y_j^2 - 2x_j y_j + x_j^2) \leq r^2$

$\iff -\left(\underbrace{\sum_{j=1}^{d} y_j^2 - 2\sum_{j=1}^{d} x_j y_j + \sum_{j=1}^{d} x_j^2 - r^2}_{f(y_1, \dots, y_d)}\right) \geq 0$

$f_1(y_1, \dots, y_d) = \sum_{j=1}^{d} y_j^2 \qquad f_2(y_1, \dots, y_d) = y_1, \qquad \dots f_{d+1}(y_1, \dots y_d) = y_d$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad f_{d+2}(y_1, \dots y_d) = 1$

$f(y_1, \dots, y_d) = -1 \cdot f_1 + 2x_1 f_2 \qquad + 2x_2 f_3 + \dots + 2x_d f_{d+1} + \left(-\sum_{j=1}^{d} x_j^2 + r^2\right) f_{d+2}$
$\Rightarrow f(y_1, \dots y_d) \in L$ and $\dim(L) = d+2$

---

$C = \{x_1, \dots x_k\}$, $G$ – concept class, $G_c = \{(g(x_1) \dots g(x_k)): g \in G\}$
Assume that $VC(G) = d$. Then $\exists \{x_1, \dots, x_d\} = C_d$ st $|G_{C_d}| = 2^d$
What if $k > d$? what can we say about $|G_c|$? (beyond the fact that $|G_c| < 2^k$)
Def (The growth function)
$\qquad \mathcal{T}_G(k) = \max |G_c| \qquad c = \{x_1, \dots x_k\}$
Lemma (Shelah – Sauer – Perles – Vapnik – Chervonenkis)
$\qquad$ Let $G$ be such that $VC(G) = d$.
$\qquad$ Then $\mathcal{T}_G(k) \leq \sum_{j=0}^{d} \binom{k}{j} \qquad \binom{k}{j} = \frac{k!}{j!(k-j)!}$

If $k \leq d$, $T_G(k) \leq \sum_{j=0}^{k} \binom{k}{j} = 2^k$

For $k > d$, $\sum_{j=0}^{d} \binom{k}{j} \leq \left(\frac{ek}{d}\right)^d$

Proof: $\binom{k}{j} = \frac{k!}{j!(k-j)!} = \frac{k(k-1)\cdots(k-j+1)}{j!}$

$= \frac{k}{d} \frac{(k-1)}{d} \cdots \frac{(k-j+1)}{d} \cdot \frac{d^j}{j!} < \left(\frac{k}{d}\right)^j \frac{d^j}{j!}$

$\underset{1 < \frac{k}{d}}{}$

$\leq \left(\frac{k}{d}\right)^d \frac{d^j}{j!}$

$\sum_{j=0}^{d} \binom{k}{j} \leq \sum_{j=0}^{d} \left(\frac{k}{d}\right)^d \frac{d^j}{j!} =$

$\leq \left(\frac{k}{d}\right)^d \sum_{j=0}^{\infty} \frac{d^j}{j!} = \left(\frac{ek}{d}\right)^d$

$\underbrace{\qquad}_{e^d}$

Exercise: $\sum_{j=0}^{d} \binom{k}{j} \geq \left(\frac{k}{d}\right)^d$

## Recap

$G$ - class of binary classifiers

$T_a(m) = \max_{c = \{x_1 \dots x_m\} \subseteq S} |G_c|$

Lemma: If $VC(G) = d < \infty$, then $T_a(m) \leq \left(\frac{me}{d}\right)^d$ for all $m > d$

Example: Let $G_1, G_2$ be two classes of binary classifiers, [Prove intersection is finite]

$VC(G_1) = d_1 < \infty$, $VC(G_2) = d_2 < \infty$

Let $C_{Gi} = \{ \{x : g(x) = +1\}, g \in G_i\}$

$i = 1, 2$

$C_{G_1} \wedge C_{G_2} = \{ C_1 \cap C_2 : C_1 \in C_{G_1}, C_2 \in C_{G_2} \}$

Let $G$ be the set of all classifiers

$g(x) = \begin{cases} 1, & x \in C \text{ for } C \in C_{G_1} \wedge C_{G_2} \\ -1, & \text{else} \end{cases}$

$\downarrow$ intersection
$\downarrow$

**Proof:** $VC(G) < \infty$

idea: If we can show that $\tau_G(m) = O(m^\nu)$ for some $\nu < \infty$,
then $VC(G) < \infty$

Fix some $\{x_1 \ldots x_m\} = M$

Consider the set $\left| M \cap \{\{ x : g(x) = +1\}, g \in C_1 \} \right| \leq \boxed{\tau_{G_1}(m)}$

illustrate:
$g_t(x) = \begin{cases} +1, x \geq t \\ -1, x < t \end{cases}$

$M = (X_1, X_2, X_3)$

$M \cap \{\{x : g(x) = +1\} g_t \in G\} = \{(X_1, X_2 X_3), (X_2, X_3) X_3, \phi\}$

$G_m = \{(+1, +1, +1), (-1, +1, +1), (-1, -1, +1), (-1, -1, -1)\}$

$M \cap \{\{X : g_1(x) = +1\}, g_1 \in G\} = M_{G_1}$

$M \cap \{\{X : g_2(x) = +1\}, g_2 \in G_2\} = M_{G_2}$

$C_1 \in M_{G_1}$   ④    $\left| C_1 \cap M_{G_2} \right| \leq \left| M_{G_2} \right| = \tau_{G_2}(m)$
      ①
         ⑤ ⑩   $\Rightarrow \left| M_{C_1} \cap M_{C_2} \right| \leq \left| M_{c_1} \right| \cdot \left| M_{c_2} \right| \leq \tau_{G_1}(m) \cdot \tau_{G_2}(m)$
                    $G_1$   $G_2$   $G_1$   $G_2$
                                    lok

$\tau_G(m) \leq \tau_{G_1}(m) \cdot \tau_{G_2}(m)$
$\leq \left(\frac{me}{d_1}\right)^{d_1} \cdot \left(\frac{me}{d_2}\right)^{d_2} = O\left(m^{d_1 + d_2}\right)$
$\Rightarrow VG(G) < \infty$

**Example:** $G_1 = \{g_t^+, t \in \mathbb{R}\}$   $VC(G_1) = 1$
$G_2 = \{g_t^-, t \in \mathbb{R}\}$   $VC(G_2) = 1$

$G = \{g_{[a,b]}, a, b \in \mathbb{R}\}$
$g_{[a,b]}(x) = \begin{cases} +1, x \in [a,b] \\ -1, else \end{cases}$

$\tau_G(m) \leq \tau_{G_1}(m) \tau_{G_2}(m)$
$\tau_{G_1}(m) = m+1, \quad \tau_{G_2}(m) = m+1$
$\tau_G : \binom{m+1}{2} = \frac{m(m+1)}{2} \leq (m+1)^2$

$\begin{array}{ccc} +1 & +1 & +1 \\ -1 & +1 & +1 \\ -1 & -1 & +1 \\ -1 & -1 & -1 \\ +1 & +1 & -1 \\ +1 & -1 & -1 \end{array}$

$(X_1, Y_1) \ldots (X_n, Y_n)$ is $\varepsilon$-representative if

$$\max_{g \in G} \left| \frac{1}{n} \sum_{1}^{n} I\{Y_j \neq g(X_j)\} - \underset{\shortparallel}{L(g)} \right| \leq \varepsilon$$

$$Pr(Y \neq g(x))$$

<u>Theorem</u>: Let $G$ be a concept class (a class of binary classifiers) and let $\tau_G(n)$ be its growth function.

Then $\max_{g \in G} \left| \frac{1}{n} \sum_{1}^{n} I\{Y_j \neq g(x_j)\} - L(g) \right| \leq \dfrac{\sqrt{2 \log(2 \tau_G(n))}}{\delta \sqrt{n}}$  $\leq \frac{\varepsilon}{?}$

with probability at least $1 - \delta$.

It's sufficient to prove that $\underset{g \in G}{\mathbb{E} \max} \left| \frac{1}{n} \sum_{s=1}^{n} I\{Y_j \neq g(x_j)\} - L(g) \right| \leq \dfrac{\sqrt{2 \log(2\tau_G(n))}}{\sqrt{n}}$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}_{Z}$

It's because

$Pr(z > t) \leq \dfrac{\mathbb{E}z}{t}$   — Markov's inequality.

$Pr(z > \underbrace{\frac{\mathbb{E}z}{\delta}}_{t}) \leq \frac{\mathbb{E}z}{\mathbb{E}z} \delta = \delta$

For (agnostic) PAC-learnability, we need

(a) An algorithm $\mathcal{A}$

(b) $n(\varepsilon, \delta)$

s.t. given $(X_1, Y_1) \ldots (X_n, Y_n)$ with $n \geq n(\varepsilon, \delta)$, $\mathcal{A}$ outputs $\hat{g}_n$ st

· $Pr(L(\hat{g}_n) > \varepsilon) \leq \delta$

In our case, if $\mathcal{A}$ is ERM, we know that $\frac{\varepsilon}{2}$-representative sample yields a classifier $L(\hat{g}_n) \leq \varepsilon$

Doing some algebra, we get $n \geq K \dfrac{V}{\delta^2 \varepsilon^2} \log\left(\dfrac{V}{\delta^2 \varepsilon^2}\right)$, $K =$ constant

Symmetrization inequality:

let $\sigma_1, ..., \sigma_n$ be random signs.
ie iid random variables st $P_r(\sigma = 1)$
$= P_r(\sigma = -1) = \frac{1}{2}$

$$\mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} I\left\{ Y_j \neq g(X_j) \right\} - L(g) \right| \leq 2 \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{1}^{n} \sigma_j I\left\{ Y_j \neq g(x_j) \right\} \right|$$

inner product $\langle (\sigma_1, ..., \sigma_n), (I\{Y_1 \neq g(X_1)\} ... I\{Y_n \neq g(x_n)\}) \rangle$

like random noise

---

Theorem  $G$ - a class of binary classifiers

$$\max_{g \in G} \left| L_n(g) - L(g) \right| \leq \frac{4}{\delta} \sqrt{\frac{\log(2 T_a(n))}{n}}$$

with probability $\geq 1 - \delta$ over the choice of the sample $(X_1, Y_1) ... (X_n, Y_n)$

In other words, if $\to$ sample size $n$ as a function of $\varepsilon, \delta$

$$n(\varepsilon, \delta) \overset{\hat{n}}{\underset{}{\geq}} 100 \, \frac{V(G)}{\delta^2 \varepsilon^2} \log\left( \frac{V}{\delta^2 \varepsilon^2} \right)$$

$\Rightarrow \hat{g}_n$ produced by ERM when given a sample of size $n(\varepsilon, \delta)$ satisfies

$$Pr(L(\hat{g}_n) > \min_{g \in G} L(g) + \varepsilon) \leq \delta \qquad \to \text{defn of PAC-learnability}$$

It suffices to show that

$$\mathbb{E} \max_{g \in G} \left| L_n(g) - L(g) \right| \leq \frac{4}{\sqrt{n}} \sqrt{\log(2 T_a(n))}$$

Symmetrization inequality:

$$\mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} I\left\{ Y_j \neq g(X_j) \right\} - L(g) \right| \leq 2 \, \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} \sigma_j I\left\{ Y_j \neq g(x_j) \right\} \right|$$

$\sigma_1, ..., \sigma_n$ iid random signs independent from $(X_1, Y_1) ... (X_n, Y_n)$, i.e. $Pr(\sigma_i = 1) \geq P_r(\sigma_i = -1)$
$= \frac{1}{2}$

Note that $\mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{1}^{n} \sigma_j I\left\{ Y_j \neq g(x_j) \right\} \right|$

$$\mathbb{E}_{(x_j, y_j)_{j=1}^n} \; \mathbb{E}_{\sigma_1 \dots \sigma_n} \; \max_{g \in G} \left| \frac{1}{n} \sum_i^n \sigma_j \, I \left\{ \overbrace{y_j \neq g(x_j)}^{t_j \in \{0, 1\}} \right\} \right|$$

<u>focus on this.</u>

$t = (t_1, \dots, t_n)$

$$\mathbb{E}_{\sigma_1 \dots \sigma_n} \; \max_{t \in T} \left| \frac{1}{n} \sum_i^n \sigma_j \, t_j \right|$$

<u>Remark</u>: $\left| \left\{ \left( I\{y_1 \neq g(x_1)\}, \dots, I\{y_n \neq g(x_n)\} \right), \; g \in G \right\} \right|$

$$= \left| G_c \right|, \quad C = \{x_1, \dots, x_n\}$$

$$\overset{\shortparallel}{\{(g(x_1), \dots, g(x_n)), \; g \in G\}}$$

$\sim 1 \qquad 1$

The number of such vectors is at most $T_G(n)$ !

<u>Lemma</u>: Let $t^{(1)}, \dots, t^{(k)} \in \mathbb{R}^n$

Then $\quad \mathbb{E} \max_{j=1 \dots k} \left| \frac{1}{n} \sum_i^n \sigma_i \, t_i^{(j)} \right| \leq 2 \max_{j=1 \dots k} \frac{\|t^{(j)}\|_2}{\sqrt{n}} \sqrt{\frac{\log(2k)}{n}}$

<u>Exercise</u> $\mathbb{E} \left| \frac{1}{n} \sum_1^n \sigma_j \, t_j \right| \leq \frac{1}{\sqrt{n}} \frac{\|t\|_2}{\sqrt{n}}$

<u>Proof</u>: Let $f(\lambda) = \mathbb{E} e^{\lambda \sigma_i} \leq e^{\frac{\lambda^2}{2}}$

Indeed, $\mathbb{E} e^{\lambda \sigma_i} = e^{\lambda \cdot 1} \frac{1}{2} + e^{\lambda(-1)} \frac{1}{2}$

$$= \frac{1}{2}(e^\lambda + e^{-\lambda}) = \frac{1}{2}\left(1 + \lambda + \frac{\lambda^2}{2} + \dots + \frac{\lambda^k}{k!} + \dots \; 1 - \lambda + \frac{\lambda^2}{2} + \dots + (-1)^k \frac{\lambda^k}{k!} + \dots\right)$$

$$= \frac{1}{2} \cdot 2 \sum_{k \geq 0} \frac{\lambda^{2k}}{(2k)!} \quad = \quad 1 + \frac{\lambda^2}{2!} + \frac{\lambda^4}{4!} + \dots$$

$$\lambda^{2k} = (\lambda^2)^k \qquad (2k)! = \underbrace{1 \cdot 2 \cdots k}_{k!} \underbrace{(k+1)(k+2) \cdots (2k)}_{\geq 2 \quad \geq 2 \quad \geq 2} \geq 2^k \cdot k!$$

$$\qquad \qquad \qquad \qquad \text{for } k \geq 1$$

$$\sum_{k \geq 0} \frac{\lambda^{2k}}{(2k)!} \leq \sum_{k \geq 0} \frac{(\lambda^2)^k}{k! \, 2^k} = \sum_{k \geq 0} \frac{\left(\frac{\lambda^2}{2}\right)^k}{k!} = e^{\frac{\lambda^2}{2}}$$

MGF (Moment Generating Function) of $\frac{\lambda}{n}\sum_i^n \sigma_j t_j$ :

$$\mathbb{E}e^{\lambda(\frac{1}{n}\sum_j^n \sigma_j t_j)} = \mathbb{E}\left(e^{\frac{\lambda}{n}\sigma_1 t_1} \cdots e^{\frac{\lambda}{n}\sigma_n t_n}\right)$$

$$= \mathbb{E}\, e^{\frac{\lambda}{n}\sigma_1 t_1} \times \cdots \times \mathbb{E}e^{\frac{\lambda}{n}\sigma_n t_n}$$

$$\le e^{\frac{\lambda^2 t_1^2}{2n^2}} \cdots e^{\frac{\lambda^2 t_n^2}{2n^2}} = e^{\frac{\lambda^2}{2n^2} \cdot \frac{\sum t_i^2}{n}} = e^{\frac{\lambda^2}{2n^2} \frac{\|t\|_2^2}{n}}$$

Next, $x \mapsto e^{\lambda x}$ is convex, i.e. $e^{\lambda(\sum_{j=1}^k \alpha_j x_j)} \le \alpha_1 e^{\lambda x_1} + \cdots + \alpha_n e^{\lambda x_n}$

$\alpha_1 \dots \alpha_k \ge 0$     In other words, $\underbrace{e^{\lambda \mathbb{E}Z} \le \mathbb{E}e^{\lambda Z}}_{\text{Jensen's inequality}}$     $Pr(X = x_j) = \alpha_j$
$\sum_j \alpha_j = 1$

$$Z = \max_{j=1\dots k}\left|\frac{1}{n}\sum_i^n \sigma_i t_i^{(j)}\right| \qquad (|a| = \max(a, -a))$$

$$e^{\lambda \mathbb{E}Z} \le \mathbb{E}e^{\lambda Z} = \mathbb{E}\max_{j=1\dots k}\left(\underbrace{e^{\frac{\lambda}{n}\sum_i^n \sigma_i t_i^{(j)}}}_{\ge 0}, \underbrace{e^{-\frac{\lambda}{n}\sum_i^n \sigma_i t_i^{(j)}}}_{\ge 0}\right)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{2k \text{ random variables}}$$

$$\mathbb{E}\sum_{j=1}^k \left(e^{\frac{\lambda}{n}\sum_i^n \sigma_i t_i^{(j)}} + e^{-\frac{\lambda}{n}\sum_i^n \sigma_i t_i^{(j)}}\right) \le 2\sum_{j=1}^k \exp\left(\frac{\lambda^2}{n^2}\frac{\|t^{(j)}\|_2^2}{2}\right)$$

$$\le 2k\exp\left(\frac{\lambda^2}{n^2}\max_{j=1\dots k}\frac{\|t^{(j)}\|_2^2}{2}\right)$$

$$e^{\lambda \mathbb{E}Z} \le 2k e^{\frac{\lambda^2}{n^2}\max_j \frac{\|t^{(j)}\|_2^2}{2}}$$

Take log : $\lambda\mathbb{E}Z \le \log(2k) + \frac{\lambda}{2n^2}\max_{1\dots k}\|t^{(j)}\|_2^2$

$\mathbb{E}Z \le \frac{\log(2k)}{\lambda} + \frac{\lambda}{2n^2}\max_{1\dots k}\|t^{(j)}\|_2^2$

Time for any $\lambda > 0$

$$h(\lambda) = \frac{\log(2k)}{\lambda} + \frac{\lambda}{2n^2}\max_{j=1\dots k}\|t^{(j)}\|_2^2$$

$$h'(\lambda) = 0 \iff \lambda_*^2 = \frac{\log(2k)}{\frac{2}{n}\max\frac{\|t^{(j)}\|_2^2}{n}}$$

$$h(\lambda_*) = \sqrt{2} \sqrt{\frac{\log (2k)}{2}} \max_{j=1,\dots,n} \frac{\| t^{(j)} \|_2}{\sqrt{n}}$$

$$\mathbb{E}_{\sigma_1,\dots,\sigma_n} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} \sigma_j \cdot I\{Y_j \neq g(X_j)\} \right| \leq \sqrt{2} \sqrt{\frac{\log(2 T_G(n))}{n}} \cdot 1$$

$$t^{(j)} = (I\{Y_1 \cdot \neq g(X_1)\}, \dots, I\{Y_n \neq g(X_n)\})$$
$$\| t^{(j)} \|_2 = \sqrt{n}$$

## Symmetrization inequality

Let $G$ be a class of binary classifiers, and $\sigma_1, \dots \sigma_n$ are
iid Random signs, i.e. $Pr(\sigma_i = 1) = Pr(\sigma_i = -1) = \frac{1}{2}$
Then $\mathbb{E} \max_{g \in G} \left| L_n(g) - L(g) \right| = \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} I\{Y_j \neq g(X_j)\} - Pr(Y \neq g(X)) \right|$

$$\leq 2 \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} \sigma_j \cdot I\{Y_j \neq g(X_j)\} \right|$$

Proof: Let $(X_1', Y_1'), \dots (X_n', Y_n')$ — an independent copy of $(X_1, Y_1), \dots (X_n, Y_n)$
Note that $L(g) = \mathbb{E} \left( \frac{1}{n} \sum_{j=1}^{n} I\{Y_j' \neq g(X_j')\} \right)$
Therefore, $\mathbb{E} \max_{g \in G} \left| L_n(g) - L(g) \right| = \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} I\{Y_j \neq g(Y_j)\} - \mathbb{E}_{(X',Y')} \frac{1}{n} \sum_{j=1}^{n} I\{Y_j' \neq g(X_j')\} \right|$

$$\leq \mathbb{E}_{(X,Y)} \max_{g \in G} \mathbb{E}_{(X',Y')} \left| \frac{1}{n} \sum_{j=1}^{n} I\{Y_j \neq g(X_j)\} - I\{Y_j' \neq g(X_j')\} \right|$$

$$| \mathbb{E} Z | \leq \mathbb{E} | Z | \qquad \leq \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} I\{Y_j \neq g(X_j)\} - I\{Y_j' \neq g(X_j')\} \right|$$
$$| \sum_i a_i | \leq \sum_i |a_i|$$
$$\max \mathbb{E} Z_i \leq \mathbb{E} \max_i Z_i$$

Note that we can "switch" $(X_j, Y_j)$ with $(X_j', Y_j')$ for any $j$
without changing the expectation.
Equivalently, for any fixed $q_1 \dots q_n \in \{+1, -1\}^n$

$$\leq \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} q_j \left( \mathbb{I}\{Y_j \neq g(X_j)\} - \mathbb{I}\{Y_j' \neq g(X_j')\} \right) \right| = \frac{1}{2^n} \sum_{(q_1 \dots q_n) \in \{-1, +1\}^n} \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} q_j \left( \mathbb{I}\{Y_j \neq g(X_j)\} \right. \right.$$

$$\left. \left. - \mathbb{I}\{Y_j' \neq g(X_j')\} \right| \right|$$

$$= \mathbb{E}_{\sigma_1 \dots \sigma_n} \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} \sigma_j \left( \mathbb{I}\{Y_j \neq g(X_j)\} - \mathbb{I}\{Y_j' \neq g(X_j')\} \right) \right|$$

$$|a - b| \leq |a| + |b| \qquad\qquad \leq 2 \, \mathbb{E} \max_{g \in G} \left| \frac{1}{n} \sum_{j=1}^{n} \sigma_j \mathbb{I}\{Y_j \neq g(X_j)\} \right|$$

## $\underline{END}$

The Fundamental Theorem of PAC Learning.

Let $G$ be a class of binary classifiers. Then the following conditions are equivalent:

(a) $G$ is agnostic PAC learnable via the ERM algorithm.

(b) $G$ is PAC learnable via the ERM algorithm.

(c) $G$ has the "uniform convergence" property: $\forall \varepsilon, \delta > 0$, $\exists n(\varepsilon, \delta)$ s.t. $\forall n \geq n(\varepsilon, \delta)$

$\max\limits_{g \in G} \left| L_n(g) - L(g) \right| \leq \varepsilon$ with probability at least $1 - \delta$.

(d) $G$ has finite VC dimension.

Proof: $(d) \Rightarrow (c)$

Moreover, we have shown that $n(\varepsilon, \delta) \leq$ constant $\dfrac{VC(G)}{\delta^2 \varepsilon^2} \log \left( \dfrac{e \, VC(G)}{\delta^2 \varepsilon^2} \right)$

$(c) \Rightarrow (a) \qquad (a) \Rightarrow (b) \qquad (b) \Rightarrow (d)$

### Remark

$$C_1 \frac{V + \log\left(\frac{1}{\delta}\right)}{\varepsilon^2} \leq n(\varepsilon, \delta) \leq C_2 \frac{V + \log\left(\frac{1}{\delta}\right)}{\varepsilon^2}$$

End of theory.

# Practical stuff

Learning beyond binary classification.
- What if there are 3 or more classes that the objects of interest should be classified into?
- What if the "label" Y takes values in $\mathbb{R}$?

- Our theory remains valid modulo minor changes.

"One vs All"        Cat vs Dog or Rabbit
                    yes ╱        ╲ no
                    Cat        Dog or Rabbit

"1 vs 1": if we have k possible labels, $\{1,...,k\}$, consider $\binom{k}{2}$ binary classification problems "X is in class i or X is in class j". Pick the label that gets most "+1" votes.

- We will talk about general "prediction" problems: predict the "response" Y based on the "predictor" X. <u>Prediction</u> is performed via some function $g \in G$.

<u>Example</u> Multi-label classification
        X is a test paper, $Y \in \{A, B, C, D, F\}$

<u>Example</u> general regression problem
        X = hours spent on social media/week,    Y = GPA $\in [1,4]$

- Need to generalization the notion of the <u>loss function</u>, denoted $\ell(y, g(x))$
    e.g. in binary classification, $\ell(y, g(x)) = \mathbb{I}\{y \neq g(x)\}$
        In multi-label classification, it can be $\ell(y, g(x)) = \mathbb{I}\{y \neq g(x)\}$
    We can also choose
    $$\ell(y, g(x)) = \begin{cases} 0, & y = g(x) \\ 1, & y \neq g(x) \text{ and } y=0 \\ 100, & y \neq g(x) \text{ and } y=1 \end{cases}$$

- The goal remaining as before; minimize $\mathbb{E}\,\ell\,(Y, g(x))$ over $g \in G$.

- _Example_ Regression problem: $Y \in \mathbb{R}$, $X \in \mathbb{R}^d$

$$\ell(y, g(x)) = (y - g(x))^2 \quad \leftarrow \text{why squared} \;(MLE)$$
$$G = \{ \langle w, x \rangle + b, \; w \in \mathbb{R}^d, \; b \in \mathbb{R} \}$$

- _Ex_ Assume that $X \in \mathbb{R}$, assume that $Y_j = \alpha x_j + \beta + \varepsilon_j$, $\alpha, \beta \in \mathbb{R}$
  and $\varepsilon_j$ is  (a) $N(0, \sigma^2)$
  
           (b) Laplace distribution with density $p(x) = \frac{1}{2} e^{-|x|}$, $x \in \mathbb{R}$
  
  $\varepsilon_1, \dots, \varepsilon_n$ are independent. Show that the MLE of $\alpha, \beta$ minimizes
  
  (a) $\frac{1}{n} \sum_{j=1}^{n} (Y_j - a x_j - b)^2$ over $a, b \in \mathbb{R}$
  
  (b) $\frac{1}{n} \sum_{1}^{n} | Y_j - a x_j - b |$

---

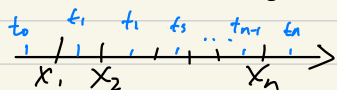Question     of practical importance :
can we implement ERM methods that have strong theoretical guarantees?

_Example_   $S = \mathbb{R}$, $T = \{ +1, -1 \}$,   $(x, y) \in S \times T$
$$G = \{ g_t^+, \; g_t^- \}, \quad g_t^+(x) = \begin{cases} +1, & x \geq t \\ -1, & x < t \end{cases}, \quad g_t^-(x) = \begin{cases} -1, & x \geq t \\ +1, & x < t \end{cases}$$

(a) _Realizable scenario_
$$\hat{g}_n \;\; \text{minimizes} \;\; \frac{1}{n} \sum_{j=1}^{n} I\{ Y_j \neq g(X_1) \} \;\; \text{over} \;\; g \in G$$



$X_{(j)}$ is the $j$th smallest among $X_1, \dots, X_n$

$O(n \log n)$ to sort. $O(\log n)$ to find $t_*$

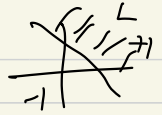Compare $L(g_{t_0}^+), \dots, L(g_{t_n}^+)$
$$L(g_{t_0}^-), \dots, L(g_{t_n}^-)$$

In agnostic learning framework, we only need to compare the empirical risks of at least $2(n+1)$ classifiers.

What about linear separators in dimension 2?
Specifically, let $S = \mathbb{R}^2$, $T = \{+1, -1\}$
$G = \{g_L, \; L \to$ a half-plane $\}$ $\quad g_L(x) = \begin{cases} 1, & x \in L \\ -1, & x \notin L \end{cases}$

Let's take a look at the more general problem:
$S = \mathbb{R}^d$, $\quad T = \{+1, -1\}$
The <u>hyperplane</u> is a set of points $\{x \in \mathbb{R}^d : \langle w, x \rangle + b = 0$
$\qquad\qquad\qquad\qquad\qquad w \in \mathbb{R}^d, \; b \in \mathbb{R} \}$

The half-spaces are given by
$L = \{ x \in \mathbb{R}^d : \langle w, x \rangle + b \geq 0$

$\tilde{x} = (x, 1)$, $\quad \tilde{w} = (w, b)$, $\quad \langle w, x \rangle + b = \langle \tilde{w}, \tilde{x} \rangle$

<u>Realizability</u> there exist $w_*$ s.t. $y_j = \text{sign}(\langle w_*, x_j \rangle)$

$\langle = \rangle \quad y_j \langle w_*, x_j \rangle > 0$

Let $\gamma = \min_{1 \dots n} y_j \langle w_*, x_j \rangle \Rightarrow y_j \langle \frac{w_*}{\gamma}, x_j \rangle \geq 1 \;\; \forall j = 1 \dots n$

$\qquad$ Denote $\tilde{w} = \frac{w_*}{\gamma} \Rightarrow y_j \langle \tilde{w}, x_j \rangle \geq 1 \;$ for all $j$.

Can we find such $\tilde{w}$?

==Perceptron Algorithm== (Frank Rosenblatt)
Given $(x_1, y_1), \dots, (x_n, y_n)$, let $w_0 = \underbrace{(0, \dots, 0)}_{d}$
for $t = 1, 2, \dots$
$\quad$ if $\exists 1 \leq j \leq n \;\; y_j \langle w^{(t)}, x_j \rangle \leq 0$
$\quad$ then $w^{(t+1)} = w^{(t)} + y_j x_j \quad$ else $\quad$ return $w^{(t)}$

$\exists w_*$ desc: $\langle w_*, x_j \rangle > 0 \iff Y_j = +1$

$$\iff Y_j \langle w_*, x_j \rangle > 0, \quad 1 \leq j \leq n$$

$$\gamma = \min_j Y_j \langle w_*, x_j \rangle \implies \min_j Y_j \langle \tfrac{w_*}{\gamma}, x_j \rangle = 1$$

Goal: find a vector $w$ s.t. $Y_j \langle w, x_j \rangle \geq 1$ for all $j$

Perceptron Algorithm (Frank Rosenblatt)

Given $(x_1, Y_1), \ldots, (x_n, Y_n)$, let $w_0 = \underbrace{(0, \ldots, 0)}_{d}$
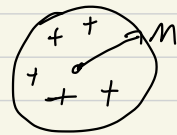
for $t = 1, 2, \ldots$

    if $\exists 1 \leq j \leq n$   $Y_j \langle w^{(t)}, x_j \rangle \leq 0$

    then $w^{(t+1)} = w^{(t)} + Y_j x_j$   else   return $w^{(t)}$

__Th__ Assume that $\max_j \| x_j \| \leq M$

Then the Perceptron algorithm stops after
at most $\left( \tfrac{M}{\gamma} \right)^2$ iterations



__Pf:__ Let $\tilde{w} = \tfrac{w_*}{\gamma}$ be s.t. $Y_j \langle \tilde{w}, x_j \rangle \geq 1$

Consider $\langle w_{t+1}, \tilde{w} \rangle = \langle w_t + Y_j x_j, \tilde{w} \rangle = \langle w_t, \tilde{w} \rangle + \underbrace{Y_j \langle x_j, \tilde{w} \rangle}_{\geq 1} \geq \langle w_t, \tilde{w} \rangle + 1$

$\implies \langle w_{t+1}, \tilde{w} \rangle \geq t + 1$

At the same time, $\| w_{t+1} \|^2 = \| w_t + Y_j x_j \|^2 = \langle w_t + Y_j x_j, w_t + Y_j x_j \rangle$

$= \| w_t \|^2 + \| x_j \|^2 + \underbrace{2 Y_j \langle w_t, w_j \rangle}_{\leq 0} \leq \| w_t \|^2 + M^2 \leq (t+1) M^2$

Combine 2 inequalities, $(t+1) \leq \langle w_{t+1}, \tilde{w} \rangle \leq \| w_{t+1} \| \cdot \| \tilde{w} \| \leq M \sqrt{t+1} \cdot \underset{\tfrac{1}{\gamma}}{\underbrace{\tfrac{1}{\gamma}}}$

$$\implies t+1 \leq \tfrac{M}{\gamma} \sqrt{t+1} \qquad \implies t+1 \leq \left( \tfrac{M}{\gamma} \right)^2$$

Perceptron as the [pseudo]-gradient descent method.

__Question:__ find $w$ s.t. $Y_j \langle w, x_j \rangle > 0$ $\forall_j$

know: $\exists \tilde{w}$ s.t. $Y_j \langle \tilde{w}, x_j \rangle \geq 1$ $\forall_j$

We can "find" $\tilde{w}$ by minimizing $F(x) = \frac{1}{2} \| x - \tilde{w} \|_2^2$ $\nabla F(x) = x - \tilde{w}$

To minimize any differentiable function $F$, we can use the gradient descent method:

let $x_0 = 0$, for $t = 1, 2, \ldots T$, $x_{t+1} = x_t - h \nabla F(x_t)$ $\nabla F(x_1, \ldots x_d) = \left( \frac{\partial F(x)}{\partial x_1}, \ldots \frac{\partial F(x)}{\partial x_d} \right)$

Pseudo-gradient descent: instead of using $\nabla F(x_t)$, assume that we can find $V_t$ such that $\langle \nabla F(x_t), V_t \rangle \geq \delta > 0$, Then define $x_{t+1} = x_t - h \cdot V_t$

$V_t = -Y_j x_j$, where $Y_j \langle w_t, x_j \rangle < 0$ for the perceptron.

$\implies \langle \nabla F(w_t), -h Y_j x_j \rangle \geq \overbrace{-h Y_j \langle w_t, x_j \rangle}^{> 0}$ $\frac{}{+ h Y_j \langle w_t, x_j \rangle} \geq h$
$\geq h$

Convergence of GD

Assume that $\forall x, y \in \mathbb{R}^d$

$\| \nabla F(x) - \nabla F(y) \| \leq L \| x - y \|_2$ E.g. if $F(x) = \frac{1}{2} \| x - \tilde{w} \|^2$, then $\nabla F(x) - \nabla F(y)$
$= x - y \implies L = 1$

Taylor's expansion:
$F(x + z) = F(x) + \langle \nabla F(\tilde{x}), z \rangle$, where $\tilde{x}$ is a point on an interval connecting $x$ and $x + z$ $\tilde{x} \overset{\bullet}{\diagup} x + z$
$\diagup^{\bullet}$

$\iff F(x+z) - F(x) = \langle \nabla F(x), z \rangle + \langle \nabla F(\tilde{x}) - \nabla F(x), z \rangle$

Let $x_{t+1} = x_t - \underbrace{h \nabla F(x_t)}_{z}$ $F(x_{t+1}) - F(x_t) = \langle \nabla F(x_t), -h \nabla F(x_t) \rangle$
$\qquad + \langle \nabla F(\tilde{x}) - \nabla F(x_t), -h \nabla F(x_t) \rangle$

$= -h \| \nabla F(x_t) \|^2 + \| \nabla F(\tilde{x}) - \nabla F(x_t) \| \cdot h \| \nabla F(x_t) \|$

$\leq -h \| \nabla F(x_t) \|^2 + L \| \tilde{x} - x_t \| \cdot h \| \nabla F(x_t) \|$

$\leq -h \| \nabla F(x_t) \|^2 + L \underbrace{\| x_{t+1} - x_t \|}_{h \| \nabla F(x_t) \|} \cdot h \| \nabla F(x_t) \|$

$\leq -h \| \nabla F(x_t) \|^2 + L h^2 \| \nabla F(x_t) \|^2$

If $h \leq \frac{1}{2L}$, then RHS $\leq -\frac{h}{2} \| \nabla F(x_t) \|^2$

$\sum_{t=0}^{\Sigma} F(x_{t+1}) - F(x_t) = F(x_{t+1}) - F(x_0) \leq \frac{-h}{2} \sum_{j=0}^{T} \| \nabla F(x_j) \|^2$

$\implies \nabla F(x_t) \to 0$ $\implies x_t \overset{t \to \infty}{\to} x_*$ s.t. $\nabla F(x_*) = 0$

$$\left\{ \; W_{t+1} = W_t + h\, Y_j X_j \right.$$

One the one hand,

$$F(W_{t+1}) = F(W_t) + \langle \nabla F(W_t), W_{t+1} - W_t \rangle + \langle \nabla F(\tilde{W}) - \nabla F(W_t), W_{t+1} - W_t \rangle$$

$$\tilde{W} \in [W_t, W_{t+1}]$$

$$\Rightarrow F(W_{t+1}) - F(W_t) \le \underbrace{h \langle \nabla F(W_t), Y_j X_j \rangle}_{\ge h} \quad \overset{h\, Y_j X_j}{}$$

$$\overset{L\, \|W_{t+1} - W_t\|_2^2}{\underset{\shortparallel}{}}$$

$$= -h + h^2 \underbrace{\|X_j\|_2^2}_{\in M} \quad \le -h + h^2 M^2$$

Take the sum for $t = 0, \dots, T$

$$F(W_{T+1}) - F(W_T) + F(W_T) - F(W_{T-1}) + \cdots$$

$$= F(W_{T+1}) - F(W_0) \le -hT + Th^2 M^2$$

Since the number of steps of perceptron does not depend on $h$

We have that

$$\frac{1}{2}\|W_{T+1} - W_*\|^2 - \frac{1}{2}\|W_*\|^2 \le -hT + Th^2 M^2$$

we know that $\|W_*\| \le \frac{1}{\gamma}$

$$hT \le Th^2 M^2 + \underset{h}{\underline{\frac{1}{2}\|W_*\|^2}} - \frac{1}{2}\|W_{T+1} - W_*\|^2$$

$$\color{blue}{\le Th M^2 + \frac{1}{2h}\|W_*\|^2 \quad \forall\, h > 0}$$

Optimize over $h \Rightarrow \sqrt{T} \le \left(\sqrt{2} + \frac{1}{\sqrt{2}}\right) \sqrt{T}\,\|W_*\| \cdot M^2$

$$T \le \left(\sqrt{2} + \frac{1}{\sqrt{2}}\right)\|W_*\|^2 \cdot M^2 = \left(\sqrt{2} + \frac{1}{\sqrt{2}}\right)\left(\frac{M}{\gamma}\right)^2$$

## Logistic Regression (an example of a "generalized linear model")

It is an example of a discriminative model: namely, it specifies the form of $P(Y | X = x)$

Here, we will assume that $X \in \mathbb{R}^d$, $Y \in \{0, 1\}$

Assume that $\Pr(Y = 1 | X = x) = p(x)$ — function of $x$

**Remark** if $p(x) > \frac{1}{2}$ $\Rightarrow$ the best guess is $Y = 1$, otherwise $Y = 0$

Note that $(X_1, Y_1)$ is the observed data, then $\mathcal{L}(p(x)) = p(x)^{Y_1}(1 - p(x_1))^{1 - Y_1}$

If the training data is $(X_1, Y_1), \ldots, (X_n, Y_n)$

then $\mathcal{L}(p(x)) = \prod_{j=1}^{n} p(x_j)^{Y_j} (1 - p(x_j))^{1 - Y_j}$

$$p(x)^Y = e^{Y \log p(x)} \quad \Longrightarrow \quad e^{\sum_{j=1}^{n} Y_j \log p(x_j) + \sum_{j=1}^{n}(1 - Y_j) \log(1 - p(x_j))}$$

$$\log \mathcal{L}(p(x)) = \sum_{j=1}^{n} Y_j \log \frac{p(x_j)}{1 - p(x_j)} + \sum_{j=1}^{n} \log(1 - p(x_j))$$

"log odds ratio"

**Main assumption:** $\log \frac{p(x)}{1 - p(x)} = \langle w, x \rangle + b = \langle \tilde{w}, \tilde{x} \rangle$

$\tilde{x} = (x, 1) \in \mathbb{R}^{d+1}$
$\tilde{w} = (w, b) \in \mathbb{R}^{d+1}$

$$\log \mathcal{L} = \sum_{j=1}^{n} Y_j \langle \tilde{w}_j, \tilde{x}_j \rangle + \sum_{j=1}^{n} \log(1 - p(x_j))$$

$$\frac{p(x)}{1 - p(x)} = e^{\langle \tilde{w}, \tilde{x} \rangle} = P(\tilde{x}) = \frac{e^{\langle \tilde{w}, \tilde{x} \rangle}}{1 + e^{\langle \tilde{w}, \tilde{x} \rangle}}$$

$$\Rightarrow 1 - p(\tilde{x}) = \frac{1}{1 + e^{\langle \tilde{w}, \tilde{x} \rangle}}$$

Therefore, maximizing $\log \mathcal{L}(p)$ is equivalent to maximizing

$$\sum_{j=1}^{n} Y_j \langle \tilde{w}_j, \tilde{x}_j \rangle - \sum_{j=1}^{n} \log\left(1 + e^{\langle \tilde{w}, \tilde{x} \rangle}\right)$$

$$\Longleftrightarrow \sum_{j=1}^{n} \log\left(1 + e^{\langle \tilde{w}_s, \tilde{x}_j \rangle}\right) = \sum_{j=1}^{n} Y_j \langle \tilde{w}_j, \tilde{x}_j \rangle \qquad \Rightarrow \text{convex}$$

minimize over $\tilde{w} \in \mathbb{R}^{d+1}$ $\Rightarrow$ it has a unique minimizer $\hat{w}$

**Remark:** a sufficient condition for $F$ to be convex is that the eigenvalues of the Hessian have to be nonnegative

Let $\hat{w}$ be the optimal solution $\Rightarrow \hat{p}(x) = \dfrac{e^{\langle \hat{w}, x \rangle}}{1 + e^{\langle \hat{w}, x \rangle}}$

$\hat{p}(x) \geqslant \frac{1}{2} \iff e^{\langle \hat{w}, x \rangle} \geqslant 1$
$\iff \langle \hat{w}, x \rangle \geqslant 0$

## Boosting

$$G = \{ g : S \to \uparrow +1, -1 \}$$

Example will be on BB

Q: What if we look at classifiers of the form
sign$[\alpha g_1 + (1 - \alpha) g_2]$, $g_1, g_2 \in G$?  $\alpha \in (0, 1)$

Recall that our goal is to find some $g$ (a binary classifier) s.t.
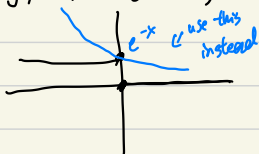$P_r(Y \neq g(x))$ is small
Note that $P_r(Y \neq g(x)) = P_r(Y g(x) < 0) = \mathbb{E} \, \mathbb{I} \{ Y g(x) < 0 \}$
Any function $f : S \to \mathbb{R}$ can be transformed into a binary classifier
$g_f : \text{sign}(f) = \begin{cases} +1, & f \geqslant 0 \\ -1, & f < 0 \end{cases}$

Problem : given a class $G$ of function $g : S \to \mathbb{R}$, minimize the empirical risk
$\frac{1}{n} \sum_{j=1}^{n} \mathbb{I} \{ Y_j \, g(x_j) < 0 \}$



Indicator function

$e^{-x}$ ← use this instead

(1) $f(x) = e^{-x}$ is convex
$f''(x) = e^{-x} > 0$

(2) $e^{-x} \geqslant \mathbb{I} \{ x \leqslant 0 \}$

Instead, consider the problem $\frac{1}{n} \sum_{j=1}^{n} e^{-Y_j \, g(x_j)} \longrightarrow$ minimize over $g \in G$
The function $z \mapsto e^{-z}$ is convex, so this can often be done numerically.

**Question**  since $\forall g \in G$, $\frac{1}{n} \sum_{j=1}^{n} e^{-Y_j \, g(x_j)} \xrightarrow{P} \mathbb{E} \, e^{-Y g(x)}$
it's natural to ask which $g$ minimizes $\mathbb{E} e^{-Y g(x)}$ over all $g : S \to \mathbb{R}$

Reminder : The minimum of $\mathbb{E}\, I\{Y_g(x) < 0\}$ is achieved for $g(x) = \text{sign}(\mathbb{E}(Y|X=x))$

**Theorem:** Let $\tilde{g}$ minimize $\mathbb{E}e^{-Yg(x)}$. Then $\text{sign}(\tilde{g}) = g*$

**Proof** : Assume that $X$ takes values $x_1, ..., x_k$. (Discrete)

$$\mathbb{E}e^{-Yg(x)} = \sum_{j=1}^{k} \mathbb{E}\left[e^{-Yg(x)} \big| X = x_k\right] Pr(X=x_k)$$

$$\mathbb{E}\left[e^{-Yg(x)} \big| X = x_k\right] = e^{1\cdot g(x_k)} P_r(Y=1|X=x_k) + e^{-1\cdot g(x_k)} P_r(Y=-1|X=x_k)$$

We know that $Pr(Y=1|X=x_k) = \dfrac{1+\eta(x_k)}{2}$   $P_r(Y=-1|X=x_k) = \dfrac{1-\eta(x_k)}{2}$

where $\eta(x_k) = \mathbb{E}[Y|X=x_k]$   $\color{blue}{\text{let } g(x_k) = t}$   $Pr(Y=t|X=x_k) = \dfrac{1+t\eta(x_k)}{2}$   $t=1$ or $-1$

Therefore, it suffices to minimize $F(t) = e^{-t}\dfrac{1+\eta(x_k)}{2} + e^{t}\dfrac{1-\eta(x_k)}{2}$ over $t \in \mathbb{R}$

$F(t) > 0$, $F''(t) = F(t) > 0$, $F$ is convex

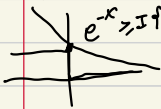$$F'(t) = -\frac{1}{e^t}\frac{1+\eta(x_k)}{2} + e^{t}\frac{1-\eta(x_k)}{2} = 0$$

$$= -(1+\eta(x_k)) + e^{2t}(1-\eta(x_k)) = 0 \implies e^{2t} = \frac{1+\eta(x_k)}{1-\eta(x_k)}, \quad t = \frac{1}{2}\log\frac{1+\eta(x_k)}{1-\eta(x_k)}$$

$F$ - "base class" of binary classifiers (e.g. threshold classifiers)

$G = \{\sum_{j=1}^{k} \alpha_j f_j, \; k \geq 1, \; \alpha_1, ..., \alpha_k \geq 0, \; f_1, ... f_k \in F\}$    $\sum_j \alpha_j = 1$

Any $g \in G$ can be transformed into a binary classifier via $g \to \text{sign}(g)$

Recall that for any binary classifier $h$, $I\{Y \neq h(x)\} = I\{Yh(x) < 0\}$


$e^{-x} \geq I\{x \leq 0\}$

$\mathbb{E}\, I\{Yg(x) < 0\}$ is minimized for $g_*(x) = \text{sign}(\mathbb{E}(Y|X=x))$

In the expression $\mathbb{E}e^{-Yg(x)}$ is minimized for $\tilde{g}(x) = \frac{1}{2}\log\frac{1+\eta(x)}{1-\eta(x)}$, where
$\eta(x) = \mathbb{E}(Y|X=x)$

$\text{sign}\left(\frac{1+\eta(x)}{1-\eta(x)}\right) = 1 \iff \frac{1+\eta(x)}{1-\eta(x)} \geq 1 \iff \eta(x) \geq 0 = \text{sign}(\eta(x))$

$\implies$ we recover the Bayes classifier!

**Summary** : minimizing $\mathbb{E}e^{-Yg(x)}$ over all functions $g$ gives us a Bayes classifier.

$\implies$ it makes sense to look at the "empirical" version of this problem,

$\frac{1}{n}\sum_{j=1}^{n} e^{-Y_j g(x_j)}$ where $(x_1, Y_1)....(x_n, Y_n)$ is the training data.

Let $G$ be a class of function, and let us consider minimizing $\frac{1}{n}\sum_i e^{-Y_i g(x_i)}$ over $g \in G$

$\color{blue}{Proof}$

<u>Definition</u> : We will say that a class $F$ of binary classifier satisfies the following

for any $n \geq 1$, any $(x_1, y_1), \ldots, (x_n, y_n)$, any non negative weights $w_1, \ldots, w_n$

s.t. $\sum_1^n w_j = 1$, $\exists f \in F$ st. $\sum_1^n w_j I \{y_j \neq f(x_j)\} \leq \frac{1}{2}$

<span style="color:blue">i.e. probability</span>

<u>Remark</u>: If $f \in F \iff -f \in F \implies$ then $F$ satisfies the weak learnability

<span style="color:blue">if loss $f > \frac{1}{2}$, we pick $-f$ which is $\leq \frac{1}{2}$ ✓ satisfied</span>

assumption.

$\frac{1}{n} \sum_1^n e^{-y_j g(x_j)} \longrightarrow$ minimize over $g \in G$

<span style="color:blue">Using the notion of pseudo gradient descent</span>

Assume that, at iteration $t$, we have $g_t \in G$

$\frac{1}{n} \sum_1^n e^{-y_i [g_t(x_j) + \alpha f(x_j)]}$

Goal: find $\alpha \in \mathbb{R}$, $f \in F$ that make this expression as small as possible

The function $f$ can be viewed as a "proxy" to the gradient.

The methods of this type are referred to as "steepest descent" methods. <span style="color:blue">umm ok.</span>

$\Rightarrow = \frac{1}{n} \sum_1^n e^{-y_j g(x_j)} e^{-y_j \alpha f(x_j)}$

Let $w_j = \frac{1}{n} e^{-y_j g(x_j)} > 0$

Then, we are trying to minimize $\sum_{j=1}^n w_j e^{-\alpha y_j f(x_j)}$ over $f \in F$, $\alpha \in \mathbb{R}$

If $\tilde{w}_j = \frac{w_j}{\sum w_j}$ so that $\sum_1^n \tilde{w} = 1$, then we need to minimize $\sum_1^n \tilde{w}_j e^{-\alpha y_j f(x_j)}$

Note that $\sum_1^n \tilde{w}_j e^{-\alpha y_j f(x_j)} = \sum_1^n \tilde{w}_j e^{-\alpha} I\{y_j = f(x_j)\} + \sum_1^n \tilde{w}_j e^{\alpha} I\{y_j \neq f(x_j)\}$

$\pm \sum \tilde{w}_j e^{-\alpha} I\{y_j \neq f(x_j)\} = e^{-\alpha} \sum \tilde{w}_j + (e^{\alpha} - e^{-\alpha}) \sum_1^n \tilde{w}_j I\{y_j \neq f(x_j)\}$

$= e^{-\alpha} + (e^{\alpha} - e^{-\alpha}) \sum_1^n \tilde{w}_j I\{y_j \neq f(x_j)\}$

To minimize this expression, (1) minimize $\sum_1^n \tilde{w}_j I\{y_j \neq f(x_j)\}$ WRT $f \in F$

(2) minimize wrt $\alpha$

Let $e_{n,t}(f) = \sum_1^n \tilde{w}_j I\{y_j \neq f(x_j)\}$

$Pr(x = x_j)$

Weak Learnability $\Rightarrow \exists \{\tilde{w}_j\}_{j=1}^n$, $\exists f \in F$ st. $e_{n,\tilde{w}}(f) \leq \frac{1}{2}$

Next, the minimum $\alpha \to e^{-\alpha} + (e^{\alpha} - e^{-\alpha}) e_{n,\tilde{w}}(f)$ is achieved

$$\hat{\alpha} = \frac{1}{2} \log \frac{1 - e_{n,\tilde{w}}(f)}{e_{n,\tilde{w}}(f)}$$

## AdaBoost

Initialize $w_j^{(1)} = \frac{1}{n}$, $j = 1, \dots, n$, $g_0 = 0$ for $t = 1, \dots, T$ do

(i) Call the weak Learner

(ii) WL outputs $f_t$ s.t. $e_{n,w^{(t)}}(f_t) \leq \frac{1}{2}$

(iii) $\alpha_t = \frac{1}{2} \log \frac{1 - e_{n,w^{(t)}}(f_t)}{e_{n,w^{(t)}}(f_t)} \geq 0$

(iv) update the weights

$$w_j^{(t+1)} = \frac{w_j^{(t)} e^{-Y_j \alpha_t f_t(x_j)}}{Z_t} \quad \leftarrow \text{normalizing factor}$$

$$Z_t = \sum_j^n w_j^{(t)} e^{-Y_j \alpha_t f_t(x_j)}$$

Output $\hat{g}_T = \frac{\sum_{t=1}^{T} \alpha_t f_t}{\sum_{t=1}^{T} \alpha_t}$

**Theorem** Assume that for any probability $w_1, \dots, w_n$, the WL finds $f$ st $\sum_j^n w_j \, \mathbb{I}\{Y_j \neq f(x_j)\} \leq \frac{1}{2} - \gamma$
for some $\gamma > 0$. Then the training error of AdaBoost satisfies
$$\frac{1}{n} \sum_j^n \mathbb{I}\{Y_j \neq \text{sign}(\hat{g}_T(x_j))\} \leq e^{-2T\gamma^2}$$

Proof: $\frac{1}{n} \sum_j^n \mathbb{I}\{Y_j \neq \text{sign}(\hat{g}_T(x_j))\} = \frac{1}{n} \sum_j^n \mathbb{I}\{Y_j \hat{g}_T(x_j) \leq 0\} \leq \frac{1}{n} \sum_j^n e^{-Y_j \hat{g}_T(x_j)} = \frac{1}{n} \sum_j^n e^{-Y_j \sum_t \alpha_t f_t(x_j)}$

Note that $w_j^{(t+1)} = \frac{1}{n} \frac{e^{Y_j \sum_i^t \alpha_i f_i(x_j)}}{\prod_i^t Z_j}$

$$e^{-Y_j \sum_i^T \alpha_i f_i(x_j)} = n \prod_i^T Z_j \, w_j^{(T+1)}$$

$$Z_t = \sum_j^n w_j^{(t)} e^{-\alpha_t Y_j f_t(x_j)}$$
$$= e^{-\alpha_t} + (e^{\alpha_t} - e^{-\alpha_t}) \sum_j^n w_j^{(t)} \mathbb{I}\{Y_j \neq f_t(x_j)\}$$

Plug in $\alpha_t = \frac{1}{2} \log \frac{1 - e_{n,w^{(t)}}(f_t)}{e_{n,w^{(t)}}(f_t)}$

$e_{n,w^{(t)}}(f_t) = \sum_j^n w_j^{(t)} \mathbb{I}\{Y_j \neq f_t(x_j)\}$

$$Z_t = 2\sqrt{e_{n,w^{(t)}}(f_t)(1 - e_{n,w^{(t)}}(f_t))}$$
$$\leq \frac{1}{2} - \gamma \leq 2\sqrt{(\frac{1}{2} - \gamma)(\frac{1}{2} + \gamma)}$$
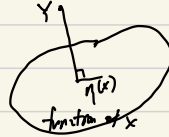
Assume that $Y$ can take values beyond $\{-1, +1\}$ (or $\{0,1\}$), specifically assume that $Y \in \mathbb{R}$

- $Y$ will be called the "response variable".
- The goal is to predict $Y$ based on the observation $X$
  - $X \in \mathbb{R}^d$, the coordinates of $X$ are called "features".
  - $X$ is also called the "predictor variable".

**Example** Predict the final exam grade $= Y$ base on (midterm grade, hw1, hw2)
$$\overbrace{\phantom{(midterm grade, hw1, hw2)}}^{\text{predictor}}$$
$$\underbrace{\phantom{grade, hw1, hw2}}_{\text{features}}$$

**Reminder :** The condition expectation of $Y$ given $X = x$, denoted $\eta(x)$,

minimize $\underset{Y|X=x}{\mathbb{E}} (Y - z)^2$

In other words, $\eta(x) = \mathbb{E}[Y | X=x]$ is the best functional approximation of $Y$ as a function of $X$



Mathematically.

$\eta(x)$ minimizes $\mathbb{E}(Y - f(x))^2$ over all functions $f$.

Given the training data $(X_1, Y_1), \ldots, (X_n, Y_n)$, we consider the problem of minimizing $\frac{1}{n} \sum_{j}^{n} (Y_j - f(X_j))^2$ over $f \in F$

**Question:** Let $\hat{f}_n$ be the solution of the problem. What is $\mathbb{E}(Y - \hat{f}_n(x))^2$ ?

Note that $\mathbb{E}(Y - \hat{f}_n(x))^2 = \mathbb{E}(Y - \eta(x) + \eta(x) - \hat{f}_n(x))^2$

orthogonal

$= \mathbb{E}(Y - \eta(x))^2 + \mathbb{E}(\eta(x) - \hat{f}_n(x))^2 + 2\underline{\mathbb{E}(Y - \eta(x))(\eta(x) - \hat{f}_n(x))}$

$= \mathbb{E}(Y - \eta(x))^2 + \mathbb{E}(\eta(x) - f_n(x))^2$

$\Downarrow$ 0

$< \mathbb{E}(Y - \eta(x))^2 + \mathbb{E}(\eta(x) - \bar{f}(x))^2 \leftarrow$ approximation error

$+ \mathbb{E}(\bar{f}(x) - f_n(x))^2 \leftarrow$ training error   $\mathbb{E}\left[\mathbb{E}[(Y - \eta(x))(\eta(x) - \hat{f}(x))|x]\right]$

$\mathbb{E}(Y|x) = \eta(x)$

Comparison to Mathematical Statistics.

$Y = \alpha X + \beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$

MLE of $(\alpha, \beta)$ is given by the solution of

$\frac{1}{n} \sum_{1}^{n} (Y_j - \alpha' X_j - \beta')^2 \longrightarrow$ minimize over $\alpha', \beta'$.

Fact: If $(Y, x)$ has bivariate normal distribution then,

$$\mathbb{E}(Y | X = x) = \alpha x + b !$$

app error

$$\mathbb{E}(\eta(x) - \hat{f}(x)) = 0$$

Error decomposition in linear regression

$$\mathbb{E}(Y - g(x))^2 = \mathbb{E}(Y - \eta(x))^2 + \underbrace{\mathbb{E}(\bar{g}(x) - \eta(x))^2}_{\text{app error of } G} + \mathbb{E}(g(x) - \bar{g}(x))^2$$

In statistics, a common assumption is that $(X, Y)$ has multivariate normal distribution. In this case, $\eta(x) = \langle w_*, x \rangle + b_*$ is a linear function of $X$, and $\underline{Y - \eta(x)}$ is normally distributed. $Y = \langle w_*, x \rangle + b_* + \varepsilon$
$\qquad \varepsilon$

· Allows to do inference: build confidence intervals / test statist hypothesis.

Solution of linear regression problem

$G = \langle g_{w,b}(x) = \langle w, x \rangle + b \rangle$

Goal. find $\hat{w}, \hat{b}$ that minimize

$\frac{1}{n} \sum_{j=1}^{n} (Y_j - \langle w, x_j \rangle - b)^2$ over $w \in \mathbb{R}^d, b \in \mathbb{R}$

Simplify: $\tilde{x}_j = (x_j, 1) \in \mathbb{R}^{d+1}$
$\qquad \tilde{w} = (w, b) \ \mathbb{R}^{d+1}$
$\qquad \langle \tilde{w}, \tilde{x}_j \rangle = \langle w, x_j \rangle + b$
$\qquad$ Let $\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \mathcal{X} = \begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \vdots \\ \tilde{x}_n \end{pmatrix}$

$$\begin{pmatrix} Y_1 - \langle \tilde{x}_1, \tilde{w} \rangle \\ Y_2 - \langle \tilde{x}_2, w \rangle \\ \vdots \\ Y_n - \langle \tilde{x}_n, w \rangle \end{pmatrix} = \vec{Y} - \mathcal{X}\tilde{w}$$

Then $\frac{1}{n}\sum_{j}^{n}(Y_j - \langle\tilde{\omega},\tilde{x}_j\rangle)^2 = \frac{1}{n}\|\vec{Y} - \mathbb{X}\tilde{\omega}\|_2^2 = F(\omega)$

If $\mathbb{X} = (x^{(1)}|x^{(2)}|\cdots|x^{(d+1)})$

$H(\omega) = \mathbb{X}\omega$ , $\nabla H(\omega) = \mathbb{X}$

$\nabla F(\omega) = -2\mathbb{X}^T(Y - \mathbb{X}\tilde{\omega}) = 0$

$\quad (\mathbb{X}^T\mathbb{X})\tilde{\omega} = X^T Y \quad$ (normal equations)

If $(\mathbb{X}^T\mathbb{X})$ is invertible

$\quad \mathbb{X} \in \mathbb{R}^{n \times (d+1)}, \ n > d+1$

$\quad \mathbb{X}^T\mathbb{X} \in \mathbb{R}^{d+1 \times d+1}$

$\therefore \hat{\omega} = (\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T Y$


Continue:

$\hat{\omega} = (X^T X)^{-1} X^T Y$ if $X^T X$ is invertible.

$\Leftrightarrow \underbrace{(X^T X)}_{A}\hat{\omega} = \underbrace{X^T Y}_{b}$

$\quad\quad A\hat{\omega} = b$

$X^T X = (X^T X)^T \Rightarrow X^T X = V\Lambda_x V^T$ where $V = (v_1|\cdots|v_p)$ is a matrix of eigen vectors.

$\Lambda_x^{-1} = \begin{pmatrix} \lambda_1^{-1} & & 0 \\ & \ddots & \\ 0 & & \lambda_p^{-1} \end{pmatrix}$ , $\lambda_1 \cdots \lambda_p$ - eigen values.

$V\Lambda_x V^T \hat{\omega} = X^T Y \quad$ Since $V^T V = I_p$, $\quad \Lambda_x V^T \hat{\omega} = V^T X^T Y \quad\quad V^T \hat{\omega} = \Lambda_x^{-1}(X^T Y)$

## The Ridge Regression

Let $\lambda > 0$ - the "regularization parameter"

$\frac{1}{n}\|\vec{Y} - \mathbb{X}\omega\|_2^2 + \lambda\|\omega\|_2^2 \longrightarrow$ minimize over $w \in \mathbb{R}^{d+1}$

$\quad\quad\quad \underbrace{\quad\quad}_{\text{regularization/penalty term}}$

$\quad\quad\quad\quad \text{Tikhonov regularization}$

$F(w) = \frac{1}{n}\|\vec{Y} - \mathbb{X}w\|_2^2 + \lambda\|w\|_2^2$

$\nabla F(w) = -\frac{2}{n}X^T(Y - Xw) + 2\lambda w$

$\nabla F(w) = 0 \Leftrightarrow \frac{2}{n}X^T Y - 2(X^T X)w + \lambda w$

$\tilde{w}$ solves the system

$$X^TXw + \lambda I \cdot w = X^TY \quad <=> \quad (X^TX + \lambda I)w = X^TY$$

$$\hat{w} = (X^TX + \lambda I)^{-1} X^TY$$

If $X^TX = V \Lambda_x V^T$, then $X^TX + \lambda I = V(\Lambda_x + \lambda I)V^T$

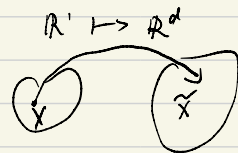Numerical Instability problem disappears, but need to pay attention to $\lambda \|w\|_2^2$.

<u>Polynomial Regression</u> (linear regression for polynomial functions).

Assume that $(x,Y) \in \mathbb{R} \times \mathbb{R}$

$$G = \{ P(x) = a_0 + a_1 x + \cdots + a_d x^d, \quad a_0 \cdots a_d \in \mathbb{R} \}$$

<u>Idea</u>: create a mapping $x \longmapsto \underbrace{(x, x^2, x^3, \ldots, x^d)}_{\tilde{x}}$

feature space

$\mathbb{R}^1 \longmapsto \mathbb{R}^d$



$$(x_1, Y_1), \cdots (x_n, Y_n) \to (\tilde{x}_1, Y_1), \ldots, (\tilde{x}_n, Y_n)$$

$$\langle w, \tilde{x} \rangle = w_1 x + w_2 x^2 + \cdots + w_d x^d$$

Linear regression problem corresponds to solving $\frac{1}{n}\sum_{i}^{n}(Y_j - \sum_{j=0}^{d} w, x^j)^2$

Also the idea of SVM

## Non-learnability of Linear Regression

Let $(X,Y) \in \mathbb{R} \times \mathbb{R}$, and $G = \{ f_w(x) = wx, \quad w \in \mathbb{R} \}$

What does it mean for $G$ to be "learnable"?

It means that $\exists$ an algorithm $A$, s.t. for any distribution over $(X,Y)$, and $\varepsilon, \delta > 0$,

$\exists n = n(\varepsilon, \delta)$, such that for all $n \geqslant n(\varepsilon, \delta)$, $A((x_1,Y_1)\ldots(x_n,Y_n))$ outputs $\hat{w}_n$ s.t.

$$\mathbb{E}(Y - \hat{w}_n x)^2 \leqslant \min_{w \in \mathbb{R}} \mathbb{E}(Y - wx)^2 + \varepsilon \quad \text{with probability} \geqslant 1 - \delta.$$

<u>Example</u> let $\varepsilon = 0.01$, $\delta = 0.5$, $n \geqslant n(\varepsilon, \delta)$

let $\mu = \dfrac{\log(\frac{100}{99})}{2n}$. Consider two distributions. $\overset{\cdot}{P_1} \quad \overset{\cdot}{P_2}$

$P_1$ : $\xrightarrow[\mu]{\overset{Y=-1}{\phantom{x}} \quad \overset{Y=0}{\phantom{x}}}_{1}$

$P_2$ : $\xrightarrow[\mu]{\overset{Y=-1}{\phantom{x}} \quad \phantom{x}}$

$\Pr((X,Y) = (1, 0)) = \mu$

$\Pr((X,Y) = (\mu, -1)) = 1$

$\Pr((X,Y) = (\mu, -1)) = 1 - \mu$

For $P_1$, $Pr(x_1 = x_2 = \dots x_n = \mu) = (1-\mu)^n \geq e^{-2\mu n} = 0.99$

Since $1-\mu \geq e^{-2\mu} = 1 - 2\mu + \frac{(2\mu)^2}{2} = 1 - 2\mu + 2\mu^2$

$1-\mu \geq 1 - 2\mu + 2\mu^2$ when $\mu \neq 0$

For $P_2$, $Pr(x_1 = \dots x_n = \mu) = 1$

We don't know whether the observation comes from which distr

$\Rightarrow A$ will produce the same output $\hat{w}_n$ regardless of the distribution.

(i) $|\hat{w}_n| < \frac{1}{2\mu}$, then $\mathbb{E}_{P_2}(Y - \hat{w}_n x)^2 = \mathbb{E}(1 - \hat{w}_n \cdot \mu)^2$ $\qquad$ $\color{blue}|\hat{w}_n \mu| < \frac{1}{2}$

$\qquad\qquad\qquad\qquad\qquad \geq (1 - \frac{1}{2})^2 = \frac{1}{4}$

But $\min_w \mathbb{E}_{P_2}(Y - wx)^2 = 0$, for $w = -\frac{1}{\mu}$

(ii) $|\hat{w}_n| \geq \frac{1}{2\mu}$ $\quad$ Consider $P_1$: $\mathbb{E}_{P_1}(Y - \hat{w}_n x)^2 = \mu(0 - \hat{w}_n \cdot 1)^2 + (1-\mu)(-1 - \hat{w}_n \mu)^2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \geq \mu \cdot w_n^2 \geq \frac{1}{4\mu}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \mu = \frac{\log(\frac{100}{99})}{2n}$

But $\min_w \mathbb{E}_{P_1}(Y - wx)^2 \leq \mathbb{E}_{P_1}(Y - 0 \cdot x)^2 = 1 - \mu$

$\begin{array}{ccc} & | & | \\ \hline & \frac{1}{1+\mu} & \end{array} \quad\quad \begin{array}{c} \\ \xrightarrow{\hspace{2cm}} \\ \frac{1}{4\mu} \end{array} \quad \Rightarrow \frac{1}{4\mu} - (1-\mu) > \varepsilon = 0.01 \quad$ for $n$ large enough

**Remark:** (a) To make the problem learnable, we need to assume that

$\qquad$ (i) $\|x\|_2 \leq M$ $\qquad$ (textbook details)

$\qquad$ (ii) $\|w\|_2 \leq R$

(b) Compare this to Gaussian linear regression:

$\qquad r = \alpha x + \varepsilon$, $x, \varepsilon$ are independent, normally distributed, then no assumption on $\alpha$ is required.

# $\mathcal{A}$rtificial $\mathcal{N}$eural $\mathcal{N}$ets

- Feedforward neural networks
  $(V, E)$ — a graph

  a set of vertices $\to$ edges
  or nodes
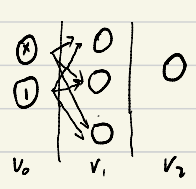
A typical graph $V = \bigcup\limits_{t=0}^{T} V_t$

$T$ - the depth of a network

Each edge in $E$ connects a vertice in $V_t$ to a vertice in $V_{t+1}$ for some $t$

Nodes correspond to "artificial neurons".

Each neuron is modeled by an "activation function" $\sigma : \mathbb{R} \to \mathbb{R}$, such as

(a) $\sigma(x) = I\{x \geq 0\}$

(b) $\sigma(x) = \dfrac{1}{1+e^{-x}}$ (sigmoid)

(c) $\sigma(x) = \max(0, x)$ (ReLU, rectified linear units).

Let $O_{t,i}(x)$ be the output of neuron $i$ in level $t$ when given input $x \in \mathbb{R}^d$

By design, $O_{0,j}(x) = x_j$, $O_{0,d+1}(x) = 1$. The input to $V_{t+1,j}$ ($j$-th neuron in layer $t+1$)

$$a_{t+1,j} = \sum W((t,r),(t+1,j)) \, O_{t,r}(x)$$
$$r(V_{t,r}, V_{t+1,j}) \in E$$

$V_0 \quad V_1 \quad V_2$

$O_{0,1}(x) = x$

$O_{0,2}(x) = 1$

$a_{1,1} = W((0,1),(1,1)) \cdot x + W((0,2),(1,1)) \cdot 1$

$a_{1,2} = W((0,1),(1,2)) \cdot x + W((0,2),(1,2)) \cdot 1$

$a_{1,3} = W((0,1),(1,3)) \cdot x + W((0,2),(1,3)) \cdot 1$

Apply activation function: $O_{1,1} = \sigma(a_{1,1})$, $O_{1,2} = \sigma(a_{1,2})$, $O_{1,3} = \sigma(a_{1,3})$

$a_{2,1} = W((1,1),(2,1)) \cdot O_{1,1} + W((1,2),(2,1)) \cdot O_{1,2} + W((1,3),(2,1)) \cdot O_{1,3}$

Explicitly: $a_{2,1} = W((1,1),(2,1)) \, \sigma(W((0,1),(1,1)) \cdot x + W((0,2),(1,1)) \cdot 1)$

$+ W((1,2),(2,1)) \, \sigma(W((0,1),(1,2)) \cdot x + W((0,2),(1,2)) \cdot 1)$

$+ W((1,3),(2,1)) \, \sigma(W((0,1),(1,3)) \cdot x + W((0,2),(1,3)) \cdot 1)$

$$G_{V,E,\sigma} = \{ g_{V,E,\sigma,w}, \quad w: E \to \mathbb{R} \} \qquad \text{Graph representation}$$
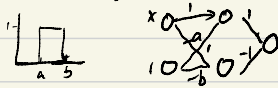
Question: how expensive can these classes be?

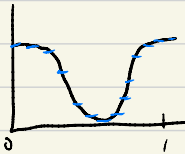Th. Let $f: [0,1] \to \mathbb{R}$ that is continuous

Then $\forall \varepsilon > 0$, $\exists (V,E)$ and weights $w \in \mathbb{R}^{|E|}$, such that $|g_{V,E,\sigma,w}(x) - f(x)| \leq \varepsilon$,

$\forall x \in [0,1]$.

We will take $\sigma(x) = I\{x \geq 0\}$.
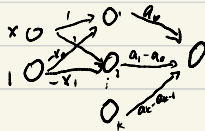
$$I\{x \in (a,b)\} = \sigma(x-a) - \sigma(x-b)$$



$$\forall \varepsilon > 0, \; \exists k \text{ and } \; 0 < x_1 < \cdots < x_k < 1 \quad \text{s.t. } \forall j \leq k, \; \left| f(x) - f\left(\frac{x_j + x_{j-1}}{2}\right) \right| \leq \varepsilon$$
$$\text{for } x \in [x_{j-1}, x_j)$$



$$\tilde{f}(x) = \sum_{j=1}^{k} a_j \, I\{x \in [x_{j-1}, x_j)\}. \quad \text{Here}, \; a_j = f\left(\frac{x_j + x_{j-1}}{2}\right)$$

Therefore, $\tilde{f}(x) = \sum_{j=1}^{k} a_j (\sigma(x - x_{j-1}) - \sigma(x - x_j))$

$$= \sum_{j=1}^{k} a_j \sigma(x - x_{j-1}) - \sum_{j=1}^{k} a_j \sigma(x - x_j)$$

$$= a_0 \sigma(x - x_0) + \sum_{j=1}^{k-1} (a_{j+1} - a_j) \sigma(x - x_j) - a_k \sigma(x - x_k)$$



1 hidden layer nn can approximate any continuous function.

final: given  , approximate it with nn
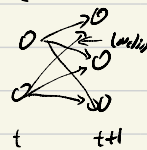
<u>Gradient Desent of NN</u>     $G = \{g_{V,E,G,W}, \; W \in \mathbb{R}^{|E|}\}$

<u>Goal</u>: $\min \; F(W) = \frac{1}{2} \sum_{j=1}^{n} (Y_j - g_W(X_j))^2$ over $W \in \mathbb{R}^{|E|}$

$(V, E), \quad V = \bigcup_{t=0}^{T} V_j \; . \quad V_t = (v_{t,1}, \dots, v_{t,k_t}), \; W_t \in \mathbb{R}^{k_{t+1} \times k_t}$

$(W_t)_{i,j}$ = height on the edge b/w $v_{t+1,i}$ and $v_{t,j}$

$W = (W_0, \dots, W_{T-1}), \; F(W) = \frac{1}{2} \sum_{j}^{n} (Y_j - g_W(X_j))^2$


$(w_t)_{i,j}$

$t \qquad t+1$

Piik some $t, \; i \le k_{t+1}, \; j \le k_t$

$\frac{\partial}{\partial (w_t)_{i,j}} F(W) = \sum_{j=1}^{n} (g_W(X_j) - Y_j) \; \frac{\partial}{\partial (w_t)_{i,j}} g_W(X_j)$

<u>Ex.</u> (a) $T=1 \implies$ no hidden layers


input    output o layer

$g_W(x) = G(W_0 \cdot X) = G(W_0 \cdot O_0)$

$\frac{\partial}{\partial (w_0)_i} g_0(x) = G'(W_0 \cdot O_0)(O_0)_i$

$\nabla_W g_W(x) = G'(W_0 \cdot O_0) O_0$

<u>Ex 2:</u>      1 hidden layer

$O_0$ - output of layer $0$

$a_1 = W_0 O_0$ - inputs of layer 1

$O_1 = G(W_0 O_0) = G(a_1) = \begin{pmatrix} G(\langle W_{0,1}, O_0 \rangle) \\ \vdots \\ G(\langle W_{0,k}, O_0 \rangle) \end{pmatrix}$

$a_2 = W_1 O_1$

$O_2 = G(a_2) = G(W_1 O_1) = G(W_1 G(W_0 O_0))$

$\nabla g_{W_1} = G'(W_1 G(W_0 O_0)) \cdot O_1$

Differentiate w.r.t. $W_0$   $\nabla g_{W_0} = G'(W_1 G(O_0 \vec{w}_0)) W_1 \, G'(O_0 \vec{w}_0) O_0$, where $G'(O_0 \vec{w}_0)$

$= \begin{pmatrix} G'(\langle W_{0,1}, O_0 \rangle) \\ \ddots \\ G'(\langle W_{0,k}, O_0 \rangle) \end{pmatrix}$

<u>Remark</u>   $W_t \in \mathbb{R}^{k_t \times k_{t-1}}, \quad k_t = \#$ of nodes in layer $t$

$\begin{pmatrix} 1 \\ 2 \\ k_t \end{pmatrix} \to \begin{pmatrix} 1 & 2 & \cdots & k_t \end{pmatrix}$

$O_{t-1} \in \mathbb{R}^{k_{t-1}} \to \begin{pmatrix} O_1^T & 0 & & 0 & 0 \\ & & & & 0 \\ 0 & O_2^T & & & 0 \\ \vdots & & \ddots & & 0 \\ 0 & & & 0 & O_{t-1}^T \end{pmatrix} = O_{t-1}$

Then $W_t O_{t-1} = O_{t-1} W_t$

In general, $\nabla_{w_0} \sigma(W_{T-1} \sigma(W_{T-2}( \dots \sigma(W_0 \cdot O_0)))$

$$= \sigma'(W_{T-1} O_{T-1}) W_{T-1} \sigma'(W_{T-2} O_{T-2}) \times \dots \times W_1 \sigma'(W_0 O_0) O_0$$