

$$D = \{(x_1, y_1), (x_2, y_2) \dots (x_m, y_m)\} \quad x \in \mathbb{R}^N \quad y \in \{-1, 1\}$$

classification: Find a function $f(x)$ such that it explains y in the best possible manner
 $\Rightarrow f(x)$ should minimize the error/loss.

argmin_f $|L(y - f(x))|$ where L stands for "loss" function. \rightarrow ①

Now, if you set something like $f(x) = \begin{cases} y_i & \text{for } i=1 \dots m \\ 0 & \text{otherwise} \end{cases}$ that will minimize the error (=0) but it is nothing but memorizing everything, and it will fail for unseen data.

Now, let's assume that x & y jointly come from a mixture distribution $F_{x,y}(x,y)$.

Obviously, since random variable x is continuous & y is discrete, marginalizⁿ w.r.t x & y will give rise to a P.D.F and a P.M.F respectively.

$$\begin{aligned} f_x &= f_x \text{ (pdf)} \\ f_y &= f_y \text{ (pmf)} \\ f_{x|y} &= f_{x|y} \text{ (pmf)} \\ f_{y|x} &= \lim_{\Delta x \rightarrow 0} \frac{P(x-\Delta x \leq x \leq x+\Delta x | y)}{P(x-\Delta x \leq x \leq x+\Delta x)} \end{aligned}$$

Notations

Prove: (Bayes theorem holds for random variables as well.)

$$f_{y|x} = \frac{f_x \cdot f_y(x|y) \cdot f_y(y)}{f_x(x)}$$

Now an ideal expansion of equⁿ (1) is to generalize by considering the expected loss function, i.e. for an unknown distribution explaining the data, the function that minimizes the expected loss (Average weighted loss for an infinite number of samples) is the best function.

$$\text{argmin}_{f(x)} = E[L(y - f(x))] = \int \int L(y - f(x)) \cdot dF_{x,y} \cdot dy \cdot dx$$

\rightarrow C.P.F (differential is a PDF).

$$\begin{aligned} &= E[E[L(y - f(x)) | f_x(x)]] \\ &= E[E[L(y - g(x)) | f_x(x)]] \\ &= \int E[L(y - g(x))] \cdot f_x(x) \cdot dx \\ &= \int \left[\int \sum_{y \in \{-1, 1\}} L(y, g(x)) \cdot f_{y|x}(y|x) \right] \cdot f_x(x) \cdot dx \end{aligned}$$

digression: $E(E(f_{y|x}(y|x)))$

$$\begin{aligned} &= \int E(f_{y|x}(y|x)) \cdot f_x(x) \cdot dx \\ &= \int \sum_{y \in \{-1, 1\}} \int f_{y|x}(y|x) \cdot f_x(x) \cdot dx \\ &= \int \sum_{y \in \{-1, 1\}} \int f_{y|x}(y|x) \cdot f_x(x) \cdot dx \\ &= \int \sum_{y \in \{-1, 1\}} f_y(y) \cdot f_x(x) \cdot dx \\ &= \int E(f_y) \cdot f_x(x) \cdot dx \\ &= E(f_y) \cdot \int f_x(x) \cdot dx \\ &= E(f_y) \cdot 1 \\ &= E(f_y) \end{aligned}$$

for binary classification:

$$\int [L(1, 1) \cdot f_{y|x}(y|x) \cdot dx + \int [L(1, -1) \cdot f_{y|x}(y|x) + L(-1, -1) \cdot f_{y|x}(y|x)]$$

The optimal solⁿ will be $g^*(x) = \begin{cases} 1 & \text{if } L(1, 1) \cdot f_{y|x}(1|x) < L(1, -1) \cdot f_{y|x}(1|x) \\ -1 & \text{else} \end{cases}$ (BAYES OPTIMAL)

Homework

Markov's inequality:

$$P(X \geq a) \leq \frac{E[X]}{a}$$

Proof: let $I(x)$ be an indicator variable $I(x) = 1$ if $x \geq a$ else 0.

$$I_{x \geq a} \equiv 1 \text{ if } x \geq a \text{ OR } 0.$$

$$\Rightarrow a I_{x \geq a} \leq X \Rightarrow E[a I_{x \geq a}] \leq E[X] \Rightarrow a E[I_{x \geq a}] \leq E[X]$$

$$\Rightarrow 1 \times P[x \geq a] + 0 \times P[x < a] \leq \frac{E[X]}{a}$$

$$\Rightarrow P[x \geq a] \leq \frac{E[X]}{a}$$

Chebyshev's inequality: $P(|x - E[x]| \geq a) \leq \frac{\text{var}(x)}{a^2}$

Strong law of large nos:

The sample average converges almost surely to the expected value.

$$\bar{x}_n \xrightarrow{\text{a.s.}} \mu \text{ when } n \rightarrow \infty \Rightarrow \text{Pr}(\lim_{n \rightarrow \infty} \bar{x}_n = \mu) = 1.$$

Weak law:

The sample average converges in probability towards the expected value.

$$\bar{x}_n \rightarrow \mu \text{ when } n \rightarrow \infty \text{ for any } \epsilon \text{ true number } \epsilon$$

$$\lim_{n \rightarrow \infty} \text{Pr}(|\bar{x}_n - \mu| > \epsilon) = 0$$

30/07/14

(LECTURE 2) dt: 30/7/14

○ We want approximate classifier that's close to the Bayesian optimal classifier.

$$E[l(y, g(x))] = \frac{1}{m} \sum_{i=1}^m l(y_i, g(x_i))$$

$$\{z_n\} \xrightarrow{P} E[X], \downarrow$$

$$S_n = P[|z_n - E[X]| > \epsilon]$$

Degenerate distribution at the end of the ~~day~~ end when the sample size is high.

$$\hat{g}_m = \min_g \frac{1}{m} \sum_{i=1}^m l(y_i, g(x_i))$$

g^* is Bayesian optimal.

(Empirical Risk min) (ERM)

g^* and \hat{g}_m^* are close in expectation,

$$E[\hat{g}_m] = g^*$$

How far is \hat{g}_m from g^* is the more imp. question.

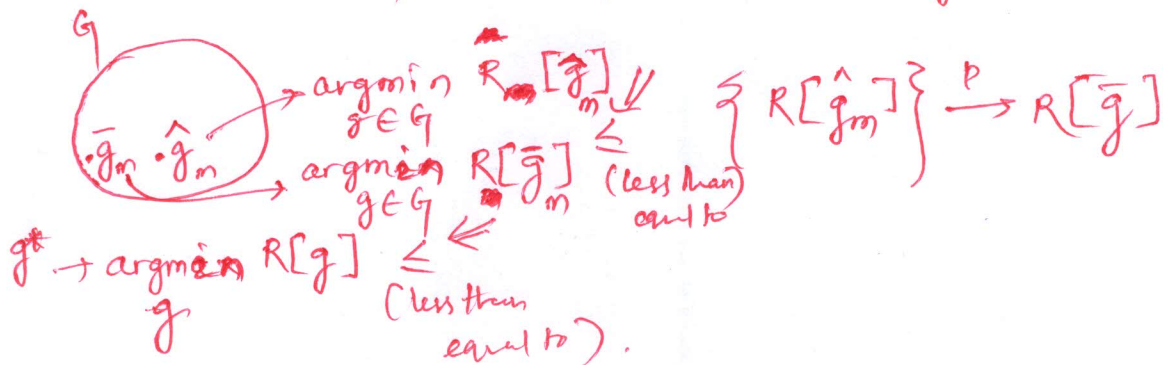
ERM is statistically constrained

$$\{R[\hat{g}_m]\} \xrightarrow{P} R[g^*]$$

True risk of empirical data should converge in optimal risk.

$\hat{R}_{10}[\hat{g}_{10}] \leq \hat{R}_{10}[g^*]$ (For a dataset as large as having 10 example, a function \hat{g}_{10} that gives the perfect solⁿ (i.e. minimizes the risk to the maximum extent) may do better than the "Bayesian optimal" (May not be the best generalization though). itself].

* So, what model we choose is more important than the algorithms.



The terms involving true risk are not easy to compute. But we have "empirical risk" which is easier to compute. Can we express in terms of empirical risk?

$$\begin{aligned}
 0 & \leq R[\hat{g}_m] - R[\bar{g}] \xrightarrow{(2)} \text{ (} \hat{R} \text{ corresponds to empirical risk)} \\
 & \leq \underbrace{[R[\hat{g}_m] - \hat{R}_m[\hat{g}_m]]}_{T} + \underbrace{[\hat{R}_m[\hat{g}_m] - R[\bar{g}]]}_{\text{2(A)}} \\
 & \quad - \hat{R}_m[\bar{g}] + \hat{R}_m[\bar{g}].
 \end{aligned}$$

(we just added & subtracted $\hat{R}_m[\hat{g}_m]$ & $\hat{R}_m[\bar{g}]$)

Since $\hat{R}_m[\bar{g}] - R[\bar{g}] \xrightarrow{P} 0$ (weak law) & $\hat{R}_m[\hat{g}_m] - \hat{R}_m[\bar{g}]$ follows equⁿ (2) [i.e. $\neq 0$], the sufficient condⁿ for (2) to be true will be if "T" is zero.

It looks like $R[\hat{g}_m] - \hat{R}_m[\hat{g}_m] \xrightarrow{P} 0$ but it may not be TRUE (why?? Homework) [Hint: Independence assumption in weak law].

Equⁿ 2(A) can be rewritten as:

$$0 \leq \max_g [R[g] - \hat{R}_m[g]] + (\text{some term} \geq 0) + 0.$$

$\hat{R}_m[g] = \Rightarrow 0$ (to be sufficient for 2).

Theorem: $\lim_{m \rightarrow \infty} P[\max_{f \in \mathcal{F}} (R[f] - \hat{R}_m[f]) > \epsilon] = 0$ if we want: $\lim_{m \rightarrow \infty} (R[f] - \hat{R}_m[f]) > \epsilon$

Weak law of large numbers:

$$\lim_{m \rightarrow \infty} P\left[\left| \frac{1}{m} \sum_{i=1}^m \ell(y_i, \hat{f}_m(x_i)) - E[\ell(y, \hat{f}_m(x))] \right| > \epsilon \right] = 0$$

$Z_i \leftarrow \bullet$

(Assuming IID)

But Z_i may not follow IID since two random variables as f may depend on all x_i 's. So Z_i 's are not independent typically.

discussion:

$$\{f_n\} \rightarrow f$$

$$\{f_n(x_i)\} \rightarrow f(x_i) \quad (\text{point wise})$$

$$\left. \begin{array}{l} |f_n(x) - f(x)| < \epsilon \quad \forall x \\ \text{OR } \max_x |f_n(x) - f(x)| < \epsilon \end{array} \right\} \text{Uniform convergence criteria.}$$

Revise

$$\text{So: } \lim_{m \rightarrow \infty} \left[\max_{f \in \mathcal{F}} (R[f] - \hat{R}_m[f]) > \epsilon \right] = 0 \rightarrow \textcircled{2}$$

(There is no "mod" so it's called one sided uniform convergence).

Illustration:

$$\mathcal{F} = \{f_1, \dots, f_n\} \quad (\text{finite}).$$

for finite set $\textcircled{2}$ can be written as:

$$\lim_{m \rightarrow \infty} P\left[\bigcup_{f \in \mathcal{F}} [R[f] - \hat{R}_m[f] > \epsilon] \right] \leq \sum_{i=1}^n P[R[f_i] - \hat{R}_m[f_i] > \epsilon]$$

$$\frac{1}{m} \sum_{j=1}^m Z_{ij} \equiv E[\ell(y, f(x))] - \ell(y_j, f_i(x_j))$$

$$\& E[Z_{ij}] = 0 \quad (\text{How})??$$

$$\sum_{i=1}^n P\left[\frac{1}{m} \sum_{j=1}^m Z_{ij} > \epsilon \right] = \sum_{i=1}^n P\left[\underbrace{\frac{1}{m} \sum_{j=1}^m Z_{ij} > \epsilon}_{(\epsilon \text{ tve.})} > \epsilon \right] \text{ OR } \sum_{i=1}^n P\left[\frac{1}{m} \sum_{j=1}^m Z_{ij} > \epsilon \right]$$

maybe -ve (can't apply Markov inequality)

$$\leq e^{-sE} \sum_{z=1}^m E \left[e^{s/m \sum_{j=1}^m z_{ij}} \right] = e^{-sE} \sum_{z=1}^m \prod_{j=1}^m E \left[e^{s/m z_{ij}} \right]$$

[Moment Generating Function of Gaussian is Gaussian.]

Fourier Transform of $\rightarrow \dots \rightarrow$]

HOFFENDING INEQ: if $E(X) = 0$, $X \in [a, b]$

$$E[e^{sX}] \leq e^{\frac{s^2(b-a)^2}{8}} \rightarrow N\left(0, \frac{(b-a)^2}{4}\right) \quad (\text{Moment generating fn of a gaussian}).$$

Nice bound.

We can also have $E[e^{sX}] \leq e^{sb}$ (But not very nice, not tighter)

$$e^{sX} \leq \frac{b-X}{b-a} e^{sa} + \frac{X-a}{b-a} e^{sb}$$

$$E(e^{sX}) \leq \left(\frac{b}{b-a}\right) e^{sa} + \left(\frac{a}{b-a}\right) e^{sb} = e^{h(z)} = -\theta z + \ln(1-\theta + \theta e^z)$$

\downarrow

$(1-\theta)e^{sa} + \theta e^{sb}$

if $z = s(b-a)$

$$= (1-\theta) \left[1 + sa + \frac{s^2 a^2}{2!} + \frac{s^3 a^3}{3!} + \dots \right] + \theta \left[e^{sb} \right]$$

$h(0) = 0$
 $h'(0) = 0$
 $h''(z) = \dots$

$$h''(z) = \frac{1}{(1-\theta + \theta e^z)^2} \times \theta e^z = \frac{1}{1-\theta + \theta e^z} \times \theta e^z + \theta e^z \times \frac{-1}{(1-\theta + \theta e^z)^2} \times \theta e^z$$

$$= \frac{\theta e^z (1-\theta)}{(1-\theta + \theta e^z)^2} \leq \frac{1}{4}$$

$h(z) \leq \frac{z^2}{8}$ (Reverse Taylor series)

$\frac{a+b}{2} \leq \sqrt{ab}$

Lecture: 4

dt: 06/08/14

Recap

① Our goal has been bringing the risk as close as possible to the optimal risk.

Statistical constraint: $\{R[\hat{f}_m]\} \xrightarrow{P} R[f]$

Bayesian constraint: $\{R[\hat{f}_m]\} \rightarrow R[f^*]$

$\lim_{m \rightarrow \infty} P[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] = 0] = 0$

(IMPORTANT POINT STARTING POINT)

② We tried taking \mathcal{F} from a finite set $\mathcal{D} \rightarrow$ i.i.d.

$$\lim_{m \rightarrow \infty} P[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] > \epsilon] = 0 \rightarrow \textcircled{1}$$

$$\lim_{m \rightarrow \infty} P \leq e^{-s \epsilon} \sum_{i=1}^m \prod_{j=1}^m E[e^{s/m Z_{ij}}]$$

(Refer to previous lecture note)
 (we applied Markov inequality)
 (we assumed i.i.d for Z_{ij})
 (we corrected the product)

$$\leq e^{-s \epsilon} \times n \times e^{s^2/m \cdot \frac{(b-a)^2}{8}} \quad \forall A > 0$$

Let's substitute $s = \frac{4m\epsilon}{(b-a)^2}$

(objective: To get tightest bound)

$$\leq n e^{-\frac{2m\epsilon^2}{(b-a)^2}} \quad (\text{Using Hoeffding inequality})$$

Comment

When $n \rightarrow \infty$, the above term will be zero. so equⁿ ① will be satisfied (∵ lower bound ≥ 0 upper bound $= 0$ (sandwiching))

$$P[R[\hat{f}_m] - \hat{R}_m[\hat{f}_m] > \epsilon] \leq n e^{-\frac{2m\epsilon^2}{(b-a)^2}} \rightarrow \textcircled{2}$$

$\rightarrow \delta$ (say) $\rightarrow \textcircled{2a}$

$$\Rightarrow P[R[\hat{f}_m] - \hat{R}_m[\hat{f}_m] \leq \epsilon] \geq 1 - \delta \rightarrow \textcircled{2(b)}$$

with probability $1 - \delta$, we have $\forall f \in \mathcal{F}$

$$R[f] \leq \hat{R}_m[f] + \epsilon \rightarrow \textcircled{3}$$

Generalizing ③ by saying that if the "max" deviation is less than ϵ then each deviation will also be less than ϵ

$$R[f] \leq \hat{R}_m[f] + (b-a) \sqrt{\frac{\log(1/\delta)}{2m}} \quad (\text{substituted } \epsilon \text{ by } \delta \text{ from } 2a)$$

$\rightarrow \textcircled{4}$

Equⁿ ④ is known as LEARNING BOUND.

(Remember it is a probabilistic bound)

The LEARNING BOUND may be satisfied for any kind of Risk minimization but we have to go back to the "starting point" and check if that holds true.

Probably Approximately Correct.

\mathcal{F} is (agnostic) PAC-learnable if for an $\epsilon > 0$; $\delta \in (0, 1)$, D , $\text{size}(\mathcal{F})$
 $x \in \mathbb{R}^D$ $m \geq \text{poly}(\frac{1}{\epsilon}, \frac{1}{\delta}, D, \text{size}(\mathcal{F})) \rightarrow P[R[\hat{f}_m] - \hat{R}_m[\hat{f}_m]] \leq \delta$

suppose we take $m = \frac{(b-a) \log \frac{1}{\delta}}{2\epsilon^2}$ (polynomial of the above described components).

Equation 2(b) will be true for our 'm' and also for $m \geq \frac{(b-a) \log \frac{1}{\delta}}{2\epsilon^2}$

So, finite \mathcal{F} is always PAC-learnable.

Example!

Let $x_1 \in (0,1)$, $x_2 \in (0,1)$, ..., $x_d \in (0,1)$ are properties of an object. and ideally if a decision is positive if $x_1 \wedge x_2 \wedge \dots \wedge x_d = f$ is 1.

so f is our f^* (optional). Now let's consider some training

example:

x_1	x_2	x_3	x_4	
1	0	1	1	\rightarrow true
0	0	1	1	\rightarrow -ve

$|f| = 3^d$ (for d terms) \rightarrow same complement sketching

what ERM can give us that we can exactly fit the training data to arrive at a "conjunction" which will not fail for training data. (or at least the loss will be minimum)

For d examples, our maximum loss can be $\frac{d}{2}$ and minimum loss can be 0. ($a=0, b=d$).

$$m \geq \frac{d}{2\epsilon^2} (d \log 3 + \log \frac{1}{\delta}) \quad (\text{PAC learnable}).$$

[PAC-learnability is a more general case and it doesn't always consider ERM]



ERM's for

Recap: A finite function classes are always statistically constraints. \rightarrow ①

$$\textcircled{1} \quad P \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] \geq \epsilon \right] \leq n e^{-2m\epsilon^2/(b-a)^2} \rightarrow \textcircled{2}$$

$$\textcircled{3} \quad m \geq \frac{(b-a)^2 \log n / \delta}{2\epsilon^2} \quad m \geq \frac{(b-a)^2 \log n / \delta}{2\epsilon^2} \rightarrow \textcircled{3} \quad [\text{Correct the previous lecture note}]$$

finite function classes are PAC learnable: \log is more helpful (why?).
 Learning Rate $\epsilon = O(1/\sqrt{m})$ [slower learning rate]

A rich function class \mathcal{F} should achieve $\hat{R}_m[f_m] = 0$

From discussion of lecture-2:

$$P \left[R[\hat{f}_m] - \hat{R}_m[\hat{f}_m] > \epsilon \right]$$

$$\leq P \left[R[\hat{f}_m] > \epsilon \right] \leq \frac{E[R[\hat{f}_m]]}{\epsilon}$$

(Using this method I can't derive easily the rate of learnability).
 (Total probability rule)

$$\leq \sum_{i: R[f_i] > \epsilon} P \left[\hat{R}_m[f_i] = 0 \right]$$

$$\leq \sum_{i: R[f_i] > \epsilon} (1 - \epsilon)^m \quad [\text{for } 0-1 \text{ loss}]$$

$$\leq n \cdot (1 - \epsilon)^m \quad (\text{will go to } 0 \text{ if } m \rightarrow \infty)$$

So statistical consistency is satisfied.

$R[f_i] > \epsilon$
 $f_1 \dots f_n$
 [For all such f_i which true risk is greater than ϵ , are candidate \hat{f}_m , so the condition $\hat{R}_m[\hat{f}_m]$ is zero.]

Using the bounding condition for any ϵ

$$1 - \epsilon \leq e^{-\epsilon} \quad (\text{use convex property of } e^{-t})$$

$$\leq n \cdot e^{-m\epsilon} \rightarrow \textcircled{4} \quad \text{learning rate } \epsilon = O\left(\frac{1}{m}\right) \quad (\text{FASTER})$$

[For non 0,1 loss, there will be some a, b terms in the express? Prove that.]

Lecture - 6

dt: 13-8-14

Arbitrary function classes (need not be finite).

$$P \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] > \epsilon \right] \leq \text{---} ?$$

In case of finite \mathcal{F} we converted to Union bounding.

$$\leq \sum_{f \in \mathcal{F}} P [R[f] - \hat{R}_m[f] > \epsilon] \quad \left(\begin{array}{l} \text{Best possible bound without} \\ \text{additional assumptions} \\ \text{Chernoff Bound} \dots \end{array} \right)$$

↳ $\frac{1}{2^m}$ equivalence classes of functions in \mathcal{F} .

Can we have a bound in case of arbitrary (non-finite) \mathcal{F} ?
What is the intuition?

→ There are many functions which will correspond to the same empirical risk (in most cases the "loss" doesn't change). The maximum of such equivalence class ~~may have~~ ^{in \mathcal{F}} ~~functions~~ will be in the order of 2^m .

$$\max_{f \in \mathcal{F}} R[f] = R[f']$$

$f' \in \mathcal{F}$

$$\therefore P \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] > \epsilon \right] \leq P \left[\max_{f' \in \mathcal{F}} R[f'] - \hat{R}_m[f] > \epsilon \right]$$

Considering the above equation and applying the notion of equivalence classes (2^m) and working out with Chernoff-bounding & Hoeffding bound we can still be able to prove that the term will converge to some value (may not be 0) irrespective of the size of \mathcal{F} (\mathcal{F} may well be ∞),

$$\textcircled{2} \quad P \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] > \epsilon \right]$$

↳ Random variable $g(u_1, \dots, u_m)$

One dimensional variation is bounded by a constant

$$u_i = (x_i, y_i)$$

↳ u_1, u_2, \dots, u_m are iid.
↳ bounded difference property.

$$|g(u_1, u_2, \dots, u_m) - g(u_1, u_2, \dots, u_k, \dots, u_m)| \leq \frac{bca}{m} \text{ (prove)}$$

(Recall)
(finite case)

$$\sum_{f \in \mathcal{F}} P[R(f) - \hat{R}_m(f)] = \frac{1}{m} \sum_{i=1}^m E(\ell(Y_i, f(x_i)) - \ell(Y_i, f(x_i)))$$

$$\sum_{i=1}^m Z_i \quad E(Z_i) = 0$$

In finite case 'g' was nothing but a sum function with variance

$$\text{var} \sum_{i=1}^m Z_i = \sum_{i=1}^m \text{var}(Z_i) \leq \left(\frac{b-a}{m}\right)^2 / 4$$

Without applying Chernoff Bound, let's go ahead..

$$P[g > \epsilon] = P[e^{sg} > e^{s\epsilon}] \leq e^{-s\epsilon} E[e^{sg}] \quad (\text{Markov})$$

$$= e^{-s\epsilon} E[E[e^{sg} | U_1, \dots, U_{m-1}]] \rightarrow \textcircled{1}$$

The introduction of conditional probability was made to enable applic of Hoeffding bound but to apply that $E[g]$ should be "zero" (which is not in our case).

Use: notation $E[g | U_1, \dots, U_{m-1}] = E(g)$

$$\max_{f \in \mathcal{F}} (\cdot) \geq (\cdot) \Rightarrow E[\max_{f \in \mathcal{F}} (\cdot)] \geq E(\cdot) \Rightarrow E[\max_{f \in \mathcal{F}} (\cdot)] \geq \max_{f \in \mathcal{F}} E(\cdot)$$

$$\rightarrow = e^{-s\epsilon} E[E[e^{s(g | U_1, \dots, U_{m-1}) - E[g | U_1, \dots, U_{m-1}]} + sE[g | U_1, \dots, U_{m-1}]]]$$

(now we can apply Hoeffding bound)

$$\leq e^{-s\epsilon} E[e^{sE[g | U_1, \dots, U_{m-1}]} \cdot e^{s^2(b-a)^2 / 8m^2}]$$

If we keep on adding and subtracting like this for $g | U_1, \dots, U_{m-2}, g | U_1, \dots, U_{m-3} \dots$ we will have $E(g)$ in the final step.

Repeated 'm' number of times, \rightarrow [Prove that these terms are also satisfying BD]

$$\leq e^{-s\epsilon} E[e^{sE[g | U_1, \dots, U_{m-1}]} \cdot e^{s^2(b-a)^2 / 8m^2}]$$

$$\leq e^{-s\epsilon} E[e^{sE[g]} \cdot e^{s^2(b-a)^2 / 8m^2 \times m}]$$

Following similar discussion as lecture-4,

$$\leq e^{-s\epsilon} \cdot e^{\frac{s^2(b-a)^2}{8m}} \cdot e^{sE(g)}$$

$$P[g > \epsilon] \leq e^{-2m[\epsilon - E(g)]^2 / (b-a)^2}$$

McDiarmid Inequality: $P[g - E(g)] \leq e^{-2\epsilon^2 / \sum c_i^2}$ & $P[E(g) - g > \epsilon] \leq e^{-2\epsilon^2 / \sum c_i^2}$

LECTURE 7

dt: 14/8/14.

Recall:
$$P \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] \right]$$

$\hookrightarrow g \rightarrow c = \frac{b-a}{m}$

McDiarmid g_i — independent rvs w/ pmf
 — bounded diff.

$$P[|g_i - E[g_i]| > \epsilon] \leq e^{-2\epsilon^2 / \sum_i c_i^2} \rightarrow \textcircled{1}$$

Refer proof of the "Bounded difference inequality" from [Berkeley.edu/wbartlett](http://Berkeley.edu/~wbartlett)

From previous lecture, we had:

$$P[g > \epsilon] \leq e^{-2m[\epsilon - E[g]]^2 / (b-a)^2} = \delta \text{ (say)}$$

We can say, with probability $1-\delta$

$$g \leq (b-a) \sqrt{\frac{\log 1/\delta}{2m}} + E[g]$$

$$\Rightarrow R[f] \leq \hat{R}_m[f] + E[g] + (b-a) \sqrt{\frac{\log 1/\delta}{2m}} \quad \forall f \in \mathcal{F}$$

$\hookrightarrow \textcircled{2}$ not finite.

(Recall) For finite case we had

$$R[f] \leq \hat{R}_m[f] + (b-a) \sqrt{\frac{\log |\mathcal{F}| / \delta}{2m}} \quad \forall f \in \mathcal{F} \text{ (finite)}$$

$\hookrightarrow \textcircled{2a}$

Speculation: This bound may be tighter than $\textcircled{2a}$ because (intuition) $\textcircled{2a}$ follows union bound (which may be unnecessary).

We will evaluate $E[g]$ later. But $E[g]$ should decay at least faster than $O(1/\sqrt{n})$ to make it better than $\textcircled{2a}$. So, choosing a function class \mathcal{F} that helps $E[g]$ decay faster is necessary.

Now let's have:

$$\Omega(f) \equiv E \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] \right] \geq 0$$

(Recall ~~max~~)

$$E(\max(\cdot)) \geq \max(E(\cdot))$$

[If: $\Omega(f_1) \geq \Omega(f_2)$ then $f_1 \succeq f_2$ (monotonic)]

~~$\Omega(f_1 \cup f_2) \geq \max(\Omega(f_1), \Omega(f_2))$~~

~~$\leq \Omega(f_1) + \Omega(f_2)$ (subadditive)]~~

$$\mathcal{R}(\mathcal{F}) = \mathbb{E} \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] \right]$$

lets say $R[f]$ is coming from some other set of data.

$$= \mathbb{E} \left[\max_{f \in \mathcal{F}} R[f] - \hat{R}_m[f] \right] \quad R_m[f] = \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i))$$

$$= \mathbb{E} \left[\max_{f \in \mathcal{F}} \mathbb{E}_{x', y'} \left[\frac{1}{m} \sum_{i=1}^m \ell(y'_i, f(x'_i)) \right] - \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i)) \right]$$

$$\leq \mathbb{E} \left[\max_{f \in \mathcal{F}} R'_m[f] - \hat{R}_m[f] \right] \rightarrow \text{discrepancy for at least one function}$$

max. discrepancy of $\mathcal{F} \Rightarrow D(\mathcal{F})$

Equⁿ (2) becomes:

$$\mathbb{E}[D(\mathcal{F})] \leq D(\mathcal{F}) + (b-a) \sqrt{\frac{\log 1/\delta}{2m}} \rightarrow (2b)$$

so with prob $1-\delta$:

$$R[f] \leq \hat{R}_m[f] + D(\mathcal{F}) + 2(b-a) \sqrt{\frac{\log 1/\delta}{2m}} \rightarrow \text{Guaranteed Risk} \rightarrow (3)$$

Intuitively,

A function class is supposed to be good if $D(\mathcal{F})$ is minimum. Even if we have at least one function in the class of \mathcal{F} which maximizes $D(\mathcal{F})$ to a large quantity, the function class will still be bad.

$$\mathbb{E} \left[\max_{f \in \mathcal{F}} \hat{R}_m[f] - \hat{R}_m[f] \right] \quad \delta = \{1, -1\} \text{ (say)}$$

$$= \mathbb{E} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \left(\sum_{i=1}^m \delta_i [\ell(y_i, f(x_i))] - \ell(y_i, f(x_i)) \right) \right]$$

$$\leq \mathbb{E} \mathbb{E}_{\delta} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta [\ell(y_i, f(x_i))] \right] + \mathbb{E} \mathbb{E}_{\delta} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m -\delta [\ell(y_i, f(x_i))] \right]$$

(Assuming δ is uniformly distributed)

$$= 2 \mathbb{E} \mathbb{E}_{\delta} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_i \ell(y_i, f(x_i)) \right]$$

(Rademacher's complexity) $\mathcal{R}_m(\mathcal{F})$

Now applying total law of Expectation on \mathcal{P}_m :

$$2 \mathbb{E}_{\mathcal{D}} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_i \ell(y_i, f(x_i)) \right]$$

$$= \mathbb{E} \left[\mathbb{E}_{\mathcal{D}} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_i \ell(y_i, f(x_i)) \mid (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m) \right] \right]$$

(Computing conditional \mathcal{P}_m is easy.)

Now applying Chernoff bounds, union bounds etc...

for a probability $1-\epsilon$:

$$\mathbb{R}[f] \leq \hat{\mathbb{R}}_m[f] + 2 \sqrt{\frac{\log 1/\epsilon}{m}} + 3(b-a) \sqrt{\frac{\log 1/\epsilon}{2m}} \quad \forall f \in \mathcal{F}$$

↳ ④

Recall:

Lecture - 8

dt: 20th Aug.

\mathcal{F} is finite

$$R[f] \leq \hat{R}_m[f] + (b-a) \sqrt{\frac{\log |\mathcal{F}|/\delta}{2m}} \quad \forall f \in \mathcal{F} \quad \rightarrow \textcircled{1}$$

$$\hookrightarrow R[f] \leq \hat{R}_m[f] + \hat{D}_m[\mathcal{F}] + (b-a) \left(\sqrt{\frac{\log 1/\delta}{2m}} + \sqrt{\frac{\log 2/\delta}{m}} \right) \quad \forall f \in \mathcal{F} \quad \rightarrow \textcircled{2}$$

$$\downarrow$$

$$\max_{f \in \mathcal{F}} \hat{R}'_m[f] - \hat{R}_m[f]$$

$$\Rightarrow R[f] \leq \hat{R}_m[f] + 2 \hat{\mathcal{R}}_m[\mathcal{F}] + 3(b-a) \sqrt{\frac{\log 2/\delta}{2m}} \quad \forall f \in \mathcal{F} \quad \rightarrow \textcircled{3}$$

$$\downarrow$$

$$E_{\sigma} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_{z_i} \ell(y_i, f(x_i)) \right] \quad (\text{Rademacher's function}).$$

$$\hat{\mathcal{R}}_m[\mathcal{F}] = E \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_i f(z_i) \right].$$

$$= E \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \delta_i (f(z_i) / \sqrt{z_i}) \right]$$

So $\hat{\mathcal{R}}$ is, as per equation $\textcircled{3}$ not only a function of function classes, it is a combination that contains the "loss" information, hence $\mathcal{R}(f)$ is sometimes replaced by $\mathcal{R}(L)$ or Rademacher function for "loss class".

Example: \mathcal{R} of linear classifiers.

lets say $L \rightarrow$ zero one loss

$$f \rightarrow \left\{ f \mid \exists w \in \mathbb{R}^N \text{ s.t. } f(x) = \text{sign}(w^T x) \quad \forall x \in \mathbb{R}^N \right\}$$

$$\hat{\mathcal{R}}(\mathcal{F}) = E \left[\max_{w \in \mathbb{R}^N} \frac{1}{m} \sum_{i=1}^m \delta_i \left(\frac{1 - y_i \text{sign}(w^T x_i)}{2} \right) \right]$$

$$\Rightarrow 2 \hat{\mathcal{R}}(\mathcal{F}) = E \left[\max_{w \in \mathbb{R}^N} \frac{1}{m} \sum_{i=1}^m \delta_i y_i \text{sign}(w^T x_i) \right] \quad \rightarrow \textcircled{4}$$

(0,1 loss func.)

Trick: Using Jensen's inequality i.e. $f(E(x)) \leq E(f(x))$,
 we can transform $E(x) \leq \frac{1}{s} \log(E[e^{sx}])$.

Equ (A) can now be transformed to:

$$2 \hat{R}(f) \leq \frac{1}{s} \ln \left(E \left[e^{s \max_{z_i} \frac{1}{m} \sum_i \delta_i z_i} \right] \right)$$

$$= \frac{1}{s} \ln \left(\max_{z_i} E \left[e^{s \frac{1}{m} \sum_i \delta_i z_i} \right] \right)$$

$$= \frac{1}{s} \ln \left(\sum_{z_i} E \left[e^{s \frac{1}{m} \sum_i \delta_i z_i} \right] \right)$$

$$= \frac{1}{s} \ln \left(\sum_{z_i} \prod_{i=1}^m E \left[e^{s/m \delta_i z_i} \right] \right)$$

$e^{\max(x)} = \max e^{(x)}$
 monotonic.

e^x is +ve
 so $\max(e^x) <$
 sum(e^x).

↳ Hoffding \rightarrow [for zero-one loss
 $\frac{b-a}{m} = 2/m$]

$$= \frac{1}{s} \ln \left(\sum_{z \in Z} e^{s^2/2m} \right) = \ln \left(\frac{|Z|}{s} \right) + \frac{s}{2m} \rightarrow \textcircled{5}$$

differencing wrt. s and equating to 0 (for finding the tightest bound),

$$2 \hat{R}(f) \leq 2 \sqrt{\frac{\log |Z|}{2m}} \rightarrow \textcircled{6}$$

$|Z|$ is referred to as growth functions of \mathcal{Z} .

(takes a value of m and returns a number)

$|Z_m|$ is upper bounded by 2^m which when replaced in $\textcircled{6}$ gives a positive value $(2 \sqrt{\log 2^m / 2})$ (not bad).

Now can we have better bound than saying $|Z_m| \leq 2^m$? (n is the dimension)

one bound is $|Z_m| \leq m C_n$ (~~n is the number of points the plane passes through~~),
 $|Z_m| \leq 2^{n+1} C_n$ (our speculation)

But according to theory, classical bound is

$$|Z_m| \leq \sum_{i=0}^{n+1} m C_i \quad (\text{prove})$$



Lecture - 9

dt: 22 Aug.

We were trying to give bound to z_m or $\Pi_f(m)$, the growth funcⁿ of linear classifiers.

$$\Pi_f(m) \begin{cases} \leq 2^{m+1} m c_n & \text{(in class)} \rightarrow \leq \left(\frac{m e}{n}\right)^n 2^{m+1} \rightarrow \textcircled{1} \\ \leq \sum_{i=0}^d m c_i & \text{(classical theory)} \rightarrow \left(\frac{m}{n+1}\right)^{n+1} \rightarrow \textcircled{2} \end{cases}$$

VC-dimensⁿ.
d is VC of \mathcal{F}

(We used $\left(\frac{m}{n}\right)^n \leq m c_n \leq \left(\frac{m e}{n}\right)^n$)

In some cases $\textcircled{1}$ may be tighter than $\textcircled{2}$ but considering all possible value of n, m & \mathcal{F} , $\textcircled{2}$ decays faster.

VC is the d where there are 'd' points that can be shattered by \mathcal{F} .

Proof of $\textcircled{2}$:

Π_f is the no. of rows

	x_1	x_2	...	x_n
f_1	+	-		+
\vdots				
f_n	-	+		-

→ distinct funcⁿ.
(at least one value should be different from other row).

↓ applying transformⁿ. (Let us assume that we can have another function class \mathcal{F}^* by changing the sign of at least one value in each row such that two rows don't become identical.)

	x_1	...	x_n
f_1	-		+
\vdots			
f_n	-	+	-

- 1 → 2 (introduce of more shattered points)
- 2 → 1 (Already retains shattered points)

By constructing examples we can say that the points which are shattered in table-1 may not be shattered in table-2 (VC dimension changes). But for those points which are shattered in 2 should be shattered in 1.

[Reason: If we change the shattered points in 1 as a result of transformⁿ, we will end up creating duplicates, so that change is not allowed, hence shattered points don't change.]

so; $\boxed{vc(\mathcal{F}) \approx vc(\mathcal{F}^*)}$. (Dimension)

we had $\hat{R}_m(\mathcal{F}) \leq \sqrt{\frac{2 \log \Pi(\mathcal{F})}{m}}$ \rightarrow PAC learnable.

so; $R[\mathcal{F}] \leq \hat{R}_m[\mathcal{F}] + 2 \sqrt{\frac{2n \log(\frac{m}{n} e)}{m}}$ (substituting 1).

Our linear classifier becomes:

$$\min_{w \in \mathbb{R}^n} \frac{1}{m} \sum_i \{ \mathbf{r}_i^T w + y_i \} \quad (\text{computationally hard to find out the combinations})$$

To make the \mathcal{F} close to Bayes optimal, problems with linear classifier

① we have to change the loss fun? (what about non-linear functions?)
if B.O is non-linear.

② To account for that if we replace x by $\phi(x)$ which will, say, ~~give~~ ^{takes} us closer to B.O. Are we doing better?

③ There is an effect of 'n' on the complexity of taking \mathcal{F} closer to Bayes optimal. "Curse of dimensionality"

[Gaussian complexity: if we replace σ_i , Rademacher's variable by a standard gaussian variable (with mean 0 & stdev. 1)]

we will have $E_{g_i} \left[\max_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m g_i f(z_i) \right]$

Binary classifiers: (Refer 0-1 loss, hinge loss, squared hinge loss, logistic loss).

[Prove $\log(1 + e^{-y f(x)})$ is convex in w .] $\left. \begin{array}{l} \text{convex} \\ \text{not diff.} \end{array} \right\}$ $\left. \begin{array}{l} \text{convex} \\ \text{diff.} \end{array} \right\}$

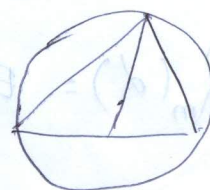
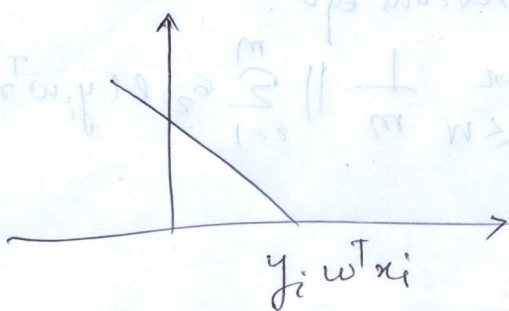
Regression: loss $\rightarrow |y - f(x_i)|$ squared loss: $(y - f(x_i))^2$

A loss function is good if it is convex, differentiable, and also address sparsity (we can ~~skip~~ ^{skip} the examples for which loss is zero if we have prior knowledge about such examples).
(y and we still get the same soln).

IMP: Most of the losses (except 0-1) are not bounded. But the data can help us put a bound on $\|w\|$ which will help us bound the loss function,
 $w^T x \leq \|w\| \|x\| \leq \|w\| R$

choice of loss func.

Hinge loss:



R - Max bound on.

V.C dimension of thick / fat classifier (for 2D case):

$$VC = \begin{cases} 1 & \text{if } M \gg 2R \\ 2 & \text{if } \frac{3R}{2} < M \leq 2R \\ 3 & \text{if } M \leq \frac{3R}{2} \end{cases} \quad \text{where } M \rightarrow \text{margin } \left(\frac{2}{\|w\|} \right)$$

when $\|w\|$ increases, $M \rightarrow 0$, we will have a large V.C dimension.

* V.C dimension can be controlled.

For 3D case:

$$VC \leq \min\left(\left\lceil \frac{4R^2}{M^2} \right\rceil, n\right) + 1 \rightarrow \text{for fat classifiers.}$$

Fat classifiers give control over V.C dimension.

(Note: Fat classifiers \neq classifiers with hinge loss).

V.C dimensions defined only for binary classes.

Radamacher's Complexity:

$$F = \left\{ f \mid \exists a \ w \in \mathbb{R}^N \rightarrow f(x) = w^T x \ \forall x, \|w\| \leq W \right\}$$

$$\hat{\Phi}_m(d) = E \left[\max_{\|w\| \leq W} \frac{1}{m} \sum_{i=1}^m \epsilon_i \max(0, 1 - y_i w^T x_i) \right]$$

Now, by Holder/Schwarz inequality:

$$\max_{\|w\| \leq W} \frac{1}{m} \sum_i \epsilon_i w^T x_i \leq W \left\| \frac{1}{m} \sum_i \epsilon_i x_i \right\|_2$$

A general form of the previous eqⁿ:

$$\hat{R}_m(\alpha) = E \left[\max_{\|w\| \leq W} \frac{1}{m} \left\| \sum_{i=1}^m \sigma_i \ell(y_i, w^T x_i) \right\|_2 \right] \rightarrow \textcircled{1}$$

Contraction lemma: Proof:

$\ell \rightarrow$ Lipschitz continuity

$$|\ell(x) - \ell(y)| \leq L \|x - y\| \quad \forall x, y,$$

(See Meir lemma: meir2003.pdf)

A convex function on any compact subset of "interior of domain" is Lipschitz continuous.

We will apply the theorem to Rademacher average of "Hinge loss" like loss function.

$$\hat{\Phi}_m(\alpha) \leq \hat{\Phi}_m(\mathcal{F}) = E \left[\max_{\|w\| \leq W} \frac{\frac{1}{m} \sum_i \delta_i y_i w^T x_i}{w^T \left(\frac{1}{m} \sum_i \delta_i y_i x_i \right)} \right] \quad (\alpha \rightarrow \text{hinge loss})$$

$$\mathcal{F} = \left\{ \text{lin functions} \mid \|w\| \leq W \right\}$$

$$\leq \frac{W}{m} E \left[\left\| \sum_i \delta_i y_i x_i \right\|^2 \right] \quad (\text{Applying Lipschitz continuity: convex lemma})$$

$$\leq \frac{W}{m} E \left[\sqrt{\left\| \sum_i \delta_i y_i x_i \right\|^2} \right]$$

Applying Jensen's inequality (since square root is convex: $E(f(x)) \leq f(E(x))$)

$$\leq \frac{W}{m} \sqrt{E \left[\left\| \sum_i \delta_i y_i x_i \right\|^2 \right]}$$

$$= \frac{W}{m} \sqrt{E \left[\sum_i \sum_j \delta_i \delta_j y_i y_j x_i x_j \right]}$$

$$= \frac{W}{m} \sqrt{\sum_i \|x_i\|^2}$$

$$\leq \frac{WR}{\sqrt{m}}$$

$$R[\mathcal{F}] \leq \hat{\Phi}_m[\mathcal{F}] + \frac{2WR}{\sqrt{m}} + 3 \sqrt{\frac{\log^2 \frac{1}{\delta}}{m}} \quad \forall f \in \mathcal{F} \rightarrow \textcircled{1}$$

Extension to non-linear classifier:

[Reason: Bayes optimal for an application need not be linear].

$$f(x) = w^T \phi(x) \quad \phi: x \rightarrow \mathbb{R}^N$$

Example: $\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2^2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$

$\phi(x)$ is not only a function giving rise to vectors with finite dimensions. $\phi(x)$ can produce a function that can take x and produce values. (may facilitate taking vectors with infinite dimension). (Reverse 'vectors', 'vector space', inner products)

(We can replace $\phi(x)$ in the above derivation (equ. 1) at least for $\phi(x)$ to be belonging to Euclidian space). More discussion to happen on this in the next lecture.

$$\min_{w \in \mathbb{R}^N} \frac{1}{m} \sum_i \ell(y_i, w^T \phi(x_j)) \quad (\text{Projected Gradient descent})$$

$$\text{s.t.: } \|w\| \leq W$$

[Terminologies: function class \equiv Model $w \equiv$ model parameters
 $W \equiv$ hyper parameter]

When we change the hyper parameters, the model changes.

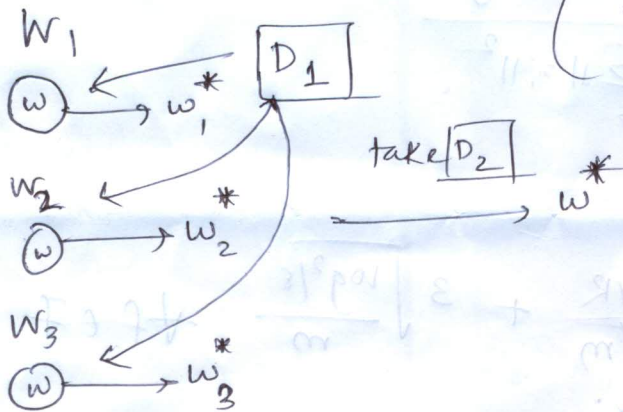
Now, how to automatically select w . \rightarrow find smallest w where $\hat{R}_m = 0$.

$$\min_{w \in \mathbb{R}^N, W} \sum_i \ell(\dots) + \frac{2WR}{\sqrt{m}}$$

$$\|w\| \leq W$$

OR
 \downarrow

Take a separate training set and do ERM for finding out w .
 (called validation).



(Surprisingly, this works better in practice).

(I) Two ways to choose W :

- ① Structural Risk Minimization: $\min_{w \in \mathbb{R}^N, w \leq W} \sum_{i=1}^N \ell(\dots) + \frac{2WR}{\sqrt{m}} \|w\| \leq W$.
- ② ERM: validation using different datasets.

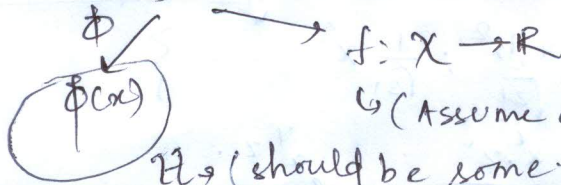
(II) ' ϕ '

- Statistical consistency $\|\phi(x)\| \leq R$ (Bounded).
- Bayes consistency ("Universal").
- X should be generic (any kind of input should be accepted).
- Computationally efficient. (A polynomial function with degree d with dimension D can grow exponentially when in terms of complexity when the input is large).

Kernels

Kernels:

Let consider X , a set (arbitrary set).



(Assume continuity)

\mathcal{H} (should be some generalization of Euclidean spaces to deal with infinite dimension).

We seek that $\phi(x)$ in the space of \mathcal{H} should be allowed to have the following properties:

$+, \cdot \rightarrow \phi(x) \subset \text{LH}(\phi(x))$ → Linear Hull (combinations),

Inner product $\langle \cdot, \cdot \rangle$ ⓐ $\langle x, x \rangle \geq 0$ ⓑ $\langle x, x \rangle = 0 \Leftrightarrow x = 0$ } Non-neg.

$\phi(x)$ itself is a function that can take parameters.
(eg $\phi(x)(z) = z$)

ⓑ symmetry: $\langle x, y \rangle = \langle y, x \rangle$

ⓒ $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$ (linearity).

$\phi(x): X \rightarrow \mathcal{H}$ (captures similarity??). → ⓐ

(The "math" way to capture similarity is to take the inner product (may be with some normalization)).

$\langle \phi(x), \phi(z) \rangle \mapsto \phi(x)(z)$ OR $f = \sum_i \alpha_i \phi(x_i)$ $g = \sum_j \beta_j \phi(x_j)$
 $\langle f, g \rangle = \sum_{i,j} \alpha_i \beta_j \phi(x_i)(x_j)$

Example: } ① $\phi(x) = x \rightarrow \phi(x)^T \phi(z) = x^T z$.
 similar representations } ② $\phi(x)(z) = x^T z \rightarrow \langle \phi(x), \phi(z) \rangle = x^T z$.
 but in the second case, ϕ is a function returning a function. Both representations are similar in terms of dot product.

③ $\phi(x)(z) = (x^T z + 1)^d$ (say), will be equivalent to $\langle \phi(x), \phi(z) \rangle$
 for $\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_1^2 \\ x_1^3 \\ \vdots \\ x_1^d \\ \sqrt{2}x_1 x_2 \end{bmatrix}$ } $\langle \phi(x), \phi(z) \rangle = (1 + x^T z)^2$ } $d=2$.

④ $\phi(x)(z) = e^{x^T z} = \langle \phi(x), \phi(z) \rangle$ for $\phi(x) = \left(\text{variant of } e^{x^T} \right)$

Proof: $\phi(x) = \left\{ 1, \frac{1}{\sqrt{2!}} x^2, \frac{1}{\sqrt{3!}} x^3, \dots \right\}$

$\langle \phi(x), \phi(z) \rangle = \sum_i \frac{x^i \cdot z^i}{i!} = e^{x^T z}$

(Remember: We are talking about the inner product in $\phi(x)$, not x).