

MACHINE LEARNING FRAMEWORK (cap. 2)

• Learner's INPUT:

- Domain set X . domain points are INSTANCES

- Label set Y .

- Training data $S \in X \times Y$. It's a sequence of labeled domain points.
(set)

• Learner's OUTPUT:

- Prediction rule $h: X \rightarrow Y$ (also called predictor or hypothesis or classifier)

$\Rightarrow A(S)$ is the hypothesis A is the learning algorithm

• Simple data-generation model:

Assume that the training data is generated by a probability distribution D over X .

(Suppose that the learner doesn't know anything about the distribution)

Suppose that there is some "correct" labeling function $F: X \rightarrow Y \Rightarrow y_i = F(x_i) \forall i$
(also unknown to the learner)

Summary: each S is generated by first sampling a point x_i according to D
and then labeling it by F .

• Measures of success: \rightarrow or RISK, GENERALIZATION ERROR, TRUE ERROR OF h

ERROR OF A CLASSIFIER: probability that it does not predict the correct label
on a random data point generated by D .

\Rightarrow the error of h is the probability to draw a random instance x , according
to D , such that $h(x) \neq F(x)$

Formally:

$A \subset X \Rightarrow D(A)$ determines how likely it is to observe a $x \in A$.

$A = \{x \in X : \pi(x) = 1\}$ with $\pi: X \rightarrow \{0, 1\}$ (we refer to A as an event)

$$D(A) \equiv P_{x \sim D} [\pi(x)]$$

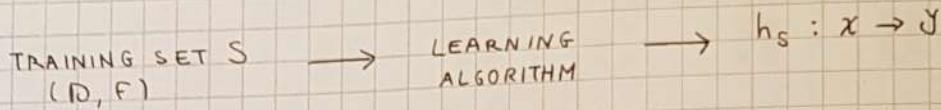
the error of $h: X \rightarrow Y$ is:

$$L_{D, F}(h) \equiv P_{x \sim D} [h(x) \neq F(x)] = D(\{x : h(x) \neq F(x)\})$$

Note that (D, F) means that the error is computed with respect to the
probability distribution D , and the correct labeling function F .

- Remainder: the learner is blind to the D and F .

EMPIRICAL RISK MINIMIZATION (ERM)



GOAL: Find h_S that minimizes the error respect to the unknown D and F .

The TRAINING ERROR is defined as: (or EMPIRICAL ERROR, EMPIRICAL RISK)

$$L_S(h) = \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m} \quad \text{where } [m] = \{1, \dots, m\}$$

Be careful w.r.t. overfitting:

EMPIRICAL RISK MINIMIZATION WITH INDUCTIVE BIAS

~~INDUCTIVE BIAS~~

Aim: find condition that guarantee no overfitting.

The learners chose before seeing the data a set of predictors $\Rightarrow H$ hypothesis class

$ERM_H(s) \in \underset{h \in H}{\operatorname{argmin}} L_S(h) \Rightarrow$ gives $h \in H$ lowest possible

In that way we "bias" the predictor toward a particular set.

• FINITE HYPOTHESIS CLASSES

H finite class

$$h_S \in \underset{h \in H}{\operatorname{argmin}} L_S(h)$$

Def (Realizability Assumption)

$$\exists h^* \in H \text{ s.t. } L_{(D,F)}(h^*) = 0 \Rightarrow L_S(h^*) = 0$$

Assumption: the points of S are generated by D independently of each other.

↑
They are INDEPENDENTLY and IDENTICALLY distributed.

i.i.d. ASSUMPTION $\rightarrow S \sim D^m$ (m size of S)

↑
probability over m -tuples to pick each element independently

S = probability of getting a non representative sample

$(1 - \delta) = \text{CONFIDENCE PARAMETER}$

$\varepsilon = \text{ACCURACY PARAMETER} \Rightarrow L_{(D,F)}(h_S) > \varepsilon, L_{(D,F)}(h_S) \leq \varepsilon$

FAILURE of the learner SUCCESS of the learner

AIM: upper bound probability to sample m-tuple of instances that will lead to failure of the learner.

upperbound: $D^m(\{S|x = (x_1, \dots, x_m) : L_{(D,F)}(h_s) > \varepsilon\}) = P_{BAD}$

set of "bad" hypothesis: $H_B = \{h \in H : L_{(D,F)}(h) > \varepsilon\}$

misleading samples: $M = \{S|x : \exists h \in H_B, L_s(h) = 0\}$

namely: $\forall S|x \in M \exists h \in H_B$ that looks like a "good" hypothesis on $S|x$.

$$\Rightarrow \{S|x : L_{(D,F)}(h_s) > \varepsilon\} \subseteq M \Rightarrow M = \bigcup_{h \in H_B} \{S|x : L_s(h) = 0\}$$

Hence:

$$D^m(\{S|x : L_{(D,F)}(h_s) > \varepsilon\}) \leq D^m(M) = D^m(\bigcup_{h \in H_B} \{S|x : L_s(h) = 0\})$$

Lemma (Union Bound)

A, B sets, D distribution: $D(A \cup B) \leq D(A) + D(B)$

Apply the UNION BOUND:

$$\Rightarrow D^m(\{S|x : L_{(D,F)}(h_s) > \varepsilon\}) \leq \sum_{h \in H_B} D^m(\{S|x : L_s(h) = 0\}) \quad (1)$$

$$D^m(\{S|x : L_s(h) = 0\}) = D^m(\{S|x : \forall i, h(x_i) = f(x_i)\}) = \prod_{i=1}^m D(\{x_i : h(x_i) = f(x_i)\})$$

For each individual sampling of an element of the training set we have

$$D(\{x_i : h(x_i) = y_i\}) = 1 - L_{(D,F)}(h) \leq 1 - \varepsilon \quad \text{in fact } h \in H_B \quad (2)$$

Combine (1) with (2) and use $1 - \varepsilon \leq e^{-\varepsilon}$:

$$\forall h \in H_B \text{ we have } D^m(\{S|x : L_s(h) = 0\}) \leq (1 - \varepsilon)^m \leq e^{-m\varepsilon} \quad (2)$$

Combining (1) with (2):

$$D^m(\{S|x : L_{(D,F)}(h_s) > \varepsilon\}) \leq |H_B| e^{-\varepsilon m} \stackrel{\text{cardinality}}{\leq} |H| e^{-\varepsilon m} \quad H_B \subset H$$

Finally:

$$P_{BAD} \leq |H| e^{-\varepsilon m} \stackrel{!}{\leq} \delta$$

$$\rightarrow e^{-\varepsilon m} \leq \frac{\delta}{|H|} \rightarrow -\varepsilon m \leq \log\left(\frac{\delta}{|H|}\right) \rightarrow m \geq -\frac{1}{\varepsilon} \log\left(\frac{\delta}{|H|}\right) \Rightarrow m \geq \frac{1}{\varepsilon} \log\left(\frac{|H|}{\delta}\right)$$

Corollary

H finite hypothesis class.

$$\delta \in (0,1), \epsilon > 0, m \geq \frac{1}{\epsilon} \log \left(\frac{|H|}{\delta} \right)$$

For any f and D , for which the realizability assumption holds (for some $h \in H$, $L(D, f)(h) = 0$), with probability of at least $1 - \delta$ over the choice of iid assumptions sample S of size m , we have that for every hypothesis h_S , it holds that :

$$L(D, f)(h_S) \leq \epsilon$$

A FORMAL LEARNING MODEL (cap. 3)

PAC LEARNING (Probably Approximately Correct (PAC))

Def (PAC Learnability)

H hypothesis class is PAC LEARNABLE if $\exists m_H : (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm with the following property

$\forall \epsilon, \delta \in (0,1), D$ distribution (over X), $f : X \rightarrow \{0,1\}$

if realizable assumption holds with respect to H, D, f

Then when running the learning algorithm on $m \geq m_H(\epsilon, \delta)$, the algorithm returns a hypothesis h such that with probability at least $1 - \delta$, $L(D, f)(h) \leq \epsilon$.

$m_H : (0,1)^2 \rightarrow \mathbb{N}$ determines the sample complexity of learning H .

→ how many examples are required to guarantee a probably approximately correct solution

Corollary

Every finite class is PAC learnable with sample complexity

$$m_H(\epsilon, \delta) \leq \left[\frac{1}{\epsilon} \log \left(\frac{|H|}{\delta} \right) \right]$$

Generalization:

- removing the realizability assumption
- learning problems beyond binary classification

AGNOSTIC PAC LEARNING

\mathbb{D} is joint probability over $X \times Y$

$\Rightarrow \mathbb{D}_X$ distribution over domain points (MARGINAL DISTRIBUTION)

$\Rightarrow \mathbb{D}((x,y) | X)$ CONDITIONAL DISTRIBUTION over labels for each domain points

• EMPIRICAL AND TRUE ERROR

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

$$L_s(h) = \frac{|\{i \in [m] : h(x_i) \neq y_i\}|}{m} \quad \text{EMPIRICAL ERROR}$$

GOAL: Find $h: X \rightarrow Y$ that minimizes the true risk $L_{\mathbb{D}}(h)$.

• BAYES OPTIMAL PREDICTOR

Probability distribution: \mathbb{D} over $X \times \{0,1\}$.

The best predictor is the BAYES OPTIMAL PREDICTOR:

$$f_{\mathbb{D}}(x) = \begin{cases} 1 & \text{if } \underbrace{\mathbb{P}[y=1 | x]}_{\text{this quantity is }} \geq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

not known

Proposition

For any classifier $g: X \rightarrow \{0,1\}$, it holds $L_{\mathbb{D}}(f_{\mathbb{D}}) \leq L_{\mathbb{D}}(g)$.

Def (AGNOSTIC PAC Learnability)

IF $\exists m_H: (0,1)^2 \rightarrow \mathbb{N}$ and a learning algorithm with the following property

$\forall \epsilon, \delta \in (0,1)$, \mathbb{D} distribution over $X \times Y$,

when running the learning algorithm on $m \geq m_H(\epsilon, \delta)$ i.i.d examples generated by \mathbb{D} , the algorithm returns h s.t., with probability $\overset{\text{of}}{\at least} 1 - \delta$

$$L_{\mathbb{D}}(h) \leq \min_{h' \in H} L_{\mathbb{D}}(h') + \epsilon$$

Idea: we drop the requirement of finding the best predictor, but we do not want to be far from it.

SCOPE OF LEARNING PROBLEMS

Domain set and training data have the same structure for all the problems.

Learner output: $h: X \rightarrow Y$ (Y is different for the 3 problems)

1) BINARY CLASSIFICATION

Target set: $Y = \{0, 1\}$

2) MULTICLASS CLASSIFICATION WITH $K > 2$ CLASSES

Target set: $Y = \{0, 1, \dots, K-1\}$

Loss: as for the binary case.

3) REGRESSION

Target set: $Y = \mathbb{R}$

Loss: need a new loss function

• GENERALIZED LOSS FUNCTION

H hypothesis class

Z domain ($x \times y$)

$\ell: H \times Z \rightarrow \mathbb{R}_+$ LOSS FUNCTION

RISK FUNCTION: expected loss of an hypothesis $h \in H$ with respect to D over Z

$$L_D(h) = E_{z \sim D} [\ell(h, z)]$$

EMPIRICAL RISK: expected loss over a given training set $S = (z_1, \dots, z_m) \in Z^m$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$$

• COMMON LOSS FUNCTION

- 0-1 loss: used in BINARY and MULTICLASS CLASSIFICATION

$$\ell_{0-1}(h, (x, y)) = \begin{cases} 0 & \text{if } h(x) = y \\ 1 & \text{if } h(x) \neq y \end{cases}$$

- square loss (L2): used in regression, penalize few large errors

$$\ell_{sq}(h, (x, y)) = (h(x) - y)^2$$

- absolute value loss (L1): used in regression, penalize many small errors

$$\ell_{abs}(h, (x, y)) = |h(x) - y|$$

Def (Ag)

H hypothesis

function

learning

H E, S

learning

D , th

1-S

L_D

where

Remark

value

LEARN

H hypothesis

ERM learn

• recu

• thi

th

Def (Agnostic PAC learnability for general loss function)

H hypothesis class, is PAC learnable with respect to a set Z and a loss function $\ell: H \times Z \rightarrow \mathbb{R}_+$, if \exists function $m_H: (0, 1)^2 \rightarrow \mathbb{N}$ and a learning algorithm with the following property

$\forall \epsilon, S \in (0, 1)$ and $\forall \mathcal{D}$ distribution over Z , when running the learning algorithm on $m \geq m_H(\epsilon, S)$ i.i.d. examples generated by \mathcal{D} , the algorithm returns $h \in H$ such that, with probability at least $1 - \delta$:

$$L_{\mathcal{D}}(h) \leq \min_{h' \in H} L_{\mathcal{D}}(h') + \epsilon$$

$$\text{where } L_{\mathcal{D}}(h) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$$

Remark $\ell(h, \cdot): Z \rightarrow \mathbb{R}_+$ random variable and $L_{\mathcal{D}}(h)$ is the expected value of this random variable.

LEARNING VIA UNIFORM CONVERGENCE (cap. 4)

H hypothesis class

ERM learning algorithm:

- receive a training sample S
- the learner evaluates the risk of each $h \in H$ on S
- the learner outputs $h^* \in H$, that minimizes this empirical risk.
- Hope: on h that minimizes the empirical risk w.r.t. S is a risk minimizer w.r.t. the true data probability distribution.

\Rightarrow SUFFICES THAT: to ensure that the empirical risk of all $h \in H$ is a good approximation of their true error

Def (ϵ -representative sample)

S training set is called ϵ -representative (w.r.t. Z, H, ℓ, \mathcal{D}) if

$$\forall h \in H, |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

Lemma

Assume: S is $\frac{\epsilon}{2}$ representative.

Then, any output of $\text{ERM}_H(S)$, namely any $h_S \in \arg\min_{h \in H} L_S(h)$, satisfies

$$L_D(h_S) \leq \min_{h \in H} L_D(h) + \epsilon$$

\Rightarrow Consequence: if with probability $1-\delta$, a random training set S is ϵ -representative, then the ERM's rule is an agnostic PAC learner.

Def (uniform convergence)

(w.r.t. Z and ϵ)

H hypothesis class has the uniform convergence property if there exists a function $m_H^{VC} : (0, 1)^2 \rightarrow \mathbb{N}$, such that $\forall \epsilon, \delta \in (0, 1)$ and $\forall D$ over Z , if S is a sample of $m \geq m_H^{VC}(\epsilon, \delta)$ examples drawn i.i.d. according to D , then with probability of at least $1-\delta$, S is ϵ -representative.

Corollary

If H has the uniform convergence property with m_H^{VC} , then the class is agnostically PAC learnable with the sample complexity $m_H(\epsilon, \delta) \leq m_H^{VC}(\frac{\epsilon}{2}, \delta)$. Furthermore, in that case the ERM_H paradigm is a successful agnostic PAC learner for H .

Corollary

Let H be a finite hypothesis class, Z a domain, $\ell : H \times Z \rightarrow [0, 1]$ be a loss function. Then:

- H enjoys the uniform convergence property with sample complexity

$$m_H^{VC}(\epsilon, \delta) \leq \left[\frac{\log(2|H|/\delta)}{2\epsilon^2} \right]$$

- H is agnostically PAC learnable using the ERM algorithm with sample complexity

$$m_H(\epsilon, \delta) \leq m_H^{VC}(\epsilon/2, \delta) \leq \left[\frac{2\log(2|H|/\delta)}{\epsilon^2} \right]$$

Proof.

$$1. \quad \mathbb{P}^m \{ S : \forall h \in H, |L_S(h) - L_D(h)| \leq \epsilon \} \geq 1-\delta$$

$$2. \quad P_{BAD} = \mathbb{P}^m \{ S : \exists h \in H, |L_S(h) - L_D(h)| > \epsilon \} \leq \delta$$

$$3. \{S : \exists h \in H, |L_S(h) - L_D(h)| > \varepsilon\} = \bigcup_{h \in H} \{S : |L_S(h) - L_D(h)| > \varepsilon\}$$

4. Union Bound:

$$\Rightarrow \Pr^m(\{S : \exists h \in H, |L_S(h) - L_D(h)| > \varepsilon\}) \leq \sum_{h \in H} \Pr^m(\{S : |L_S(h) - L_D(h)| > \varepsilon\})$$

5. Demonstrate that for any fixed h , the difference $|L_S(h) - L_D(h)|$ is small enough.

$$6. \text{Recall that: } L_D(h) = \mathbb{E}_{z \sim D} [l(h, z)] \text{ and } L_S(h) = \frac{1}{m} \sum_{i=1}^m l(h, z_i).$$

By the linearity of expectation, it follows that $L_D(h)$ is also the expected value of $L_S(h)$. Hence, the quantity $|L_D(h) - L_S(h)|$ is the deviation of the random variable $L_S(h)$ from its expectation.

→ Need to show that the measure of $L_S(h)$ is concentrated around its expected value.

7. Law of large numbers: if m is large, the average converges to the expectation.

8. Lemma (Hoeffding's inequality)

$\theta_1, \dots, \theta_m$ is a sequence of i.i.d. random variables.

Assume that $\forall i, \mathbb{E}[\theta_i] = \mu$ and $\mathbb{P}[a \leq \theta_i \leq b] = 1$.

Then for any $\varepsilon > 0$,

$$\Pr\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \varepsilon\right] \leq 2 \exp\left(-\frac{2m\varepsilon^2}{(b-a)^2}\right) \quad [0, 1]$$

$$9. \Rightarrow \Pr^m(\{S : |L_S(h) - L_D(h)| > \varepsilon\}) = \Pr\left[\left|\frac{1}{m} \sum_{i=1}^m \theta_i - \mu\right| > \varepsilon\right] \leq 2 \exp(-2m\varepsilon^2)$$

$$10. \sum_{h \in H} \Pr^m(\{S : |L_S(h) - L_D(h)| > \varepsilon\}) \leq |\mathcal{H}| 2 \exp(-2m\varepsilon^2)$$

$$= \sum_{h \in H} 2 \exp(-2m\varepsilon^2)$$

11. Force:

$$|\mathcal{H}| 2 e^{-2m\varepsilon^2} \stackrel{!}{\leq} \delta \Rightarrow m \geq \log\left(\frac{2|\mathcal{H}|}{\delta}\right) \frac{1}{2\varepsilon^2}$$

$$\Rightarrow \sum_{h \in H} \Pr^m(\{S : |L_S(h) - L_D(h)| > \varepsilon\}) \leq \delta$$

□

Remark (The discretization trick)

simple trick that allows us to get a good estimate of the practical sample complexity of infinite hypothesis class

- Consider an hypothesis class determined by d real number parameters.
- In principle: hypothesis class of infinite size
- In practice: real numbers represented with 64 bits double precision variables.

• For d parameters: $|H| = 2^{64d}$ (so $|H|$ is larger but finite)

$$m_H(\varepsilon, S) \leq m_H^{\text{VC}}\left(\frac{\varepsilon}{2}, S\right) \leq \frac{2 \log(2^{64d}/\varepsilon)}{\varepsilon^2}$$

- \Rightarrow the bound depends on the chosen number representation.

THE BIAS COMPLEXITY TRADE-OFF (cap. 5)

A learning algorithm

S training set

H hypothesis set $\Rightarrow \hat{h}$ from H for which $L_D(\hat{h})$ is small.

Question: Is there a universal learner, i.e. an algorithm A, that predicts

the best \hat{h} for any distribution D ?

NO FREE-LUNCH THEOREM

Theorem (NO FREE-LUNCH)

Let A be any learning algorithm for the task of binary classification with respect to the 0-1 loss over a domain X. Let m be any number smaller than $\frac{|X|}{2}$, representing a training set size. Then, there exists a distribution D over $X \times \{0,1\}$ such that

1. There exists a function $f: X \rightarrow \{0,1\}$ with $L_D(f) = 0$.

2. With probability of at least $\frac{1}{7}$ over the choice of $S \sim D^m$ we have that $L_D(A(s)) \geq 1/8$.

\Rightarrow meaning: there is no universal learner. In fact, no learner can succeed on all learning tasks. For every learner, there exist a task on which it fails, even though that task can be successfully learned by another learner.

Corollary

Let X be an infinite domain set and let H be the set of all functions from X to $\{0,1\}$. Then, H is not PAC learnable.

Idea of the proof of the theorem: our training set is smaller than half of the domain, so we have no information on what happens on the other half.

There exist some target function f that works on the other half in a way that contradicts our estimated labels.

Remark Class H of all possible functions from X to $\{0,1\}$ (i.e. assuming no prior knowledge) is not a good idea. \rightarrow it is not PAC LEARNABLE.

Proof of the corollary.

proceed by contradiction:

1. Assume PAC learnable.

- choose $\epsilon < 1/8$, $\delta < 1/7$ (NB. H includes all functions)
- By definition of PAC: \exists an algorithm A s.t. For every distribution D , if realizable, when running the algorithm on m i.i.d. examples generated by D , the algorithm returns a hypothesis h such that, with probability $\geq 1-\delta$, $L_D(A(s)) \leq \epsilon$.

2. Apply no free-lunch theorem to A .

- $|X| > 2m$: For any ML algorithm (including A) there exist a distribution D for which with probability $\geq 1/7 > \delta$,
 $L_D(A(s)) \geq 1/8 > \epsilon$.

3. Contradiction!

Remark (choose a good hypothesis set)

- We need to use our prior knowledge about D , to pick a good hypothesis set.
- We would like H to be large, so that it may contain a function h with small $L_S(h)$ and hopefully a small $L_D(h)$.

• No Free-lunch: H cannot be too large!

SHALL: large L_s and $L_D \approx L_s$

H
LARGE: small L_s and $L_D \gg L_s$
(risk of overfitting)

ERROR DECOMPOSITION

h_S ERM hypothesis

$$\Rightarrow L_D(h_S) = E_{app} + E_{est}$$

where

$$E_{app} = \text{approximation error} = \min_{h \in H} L_D(h)$$

Minimum risk achievable by a predictor in the hypothesis class.

Measures how much risk we have because we restrict ourselves to a specific class, namely, how much INDUCTIVE BIAS we have.

- depends on the choice of H .
- does not depend on the sample size m .
- under the realizability assumption $\Rightarrow E_{app} = 0$!
- larger $H \Rightarrow$ smaller E_{app} .

$$E_{est} = \text{estimation error} = L_D(h_S) - E_{app}$$

Difference between the approximation error and the error achieved by the ERM predictor. The estimation error results

because the empirical risk (training error) is only an estimate of the true risk, and so the predictor minimizing the empirical risk is only an estimate of the predictor minimizing the true risk.

- depends on the training set size m and on the size (complexity) of the hypothesis class;

E_{est} increase with $|H|$ and decrease with m .

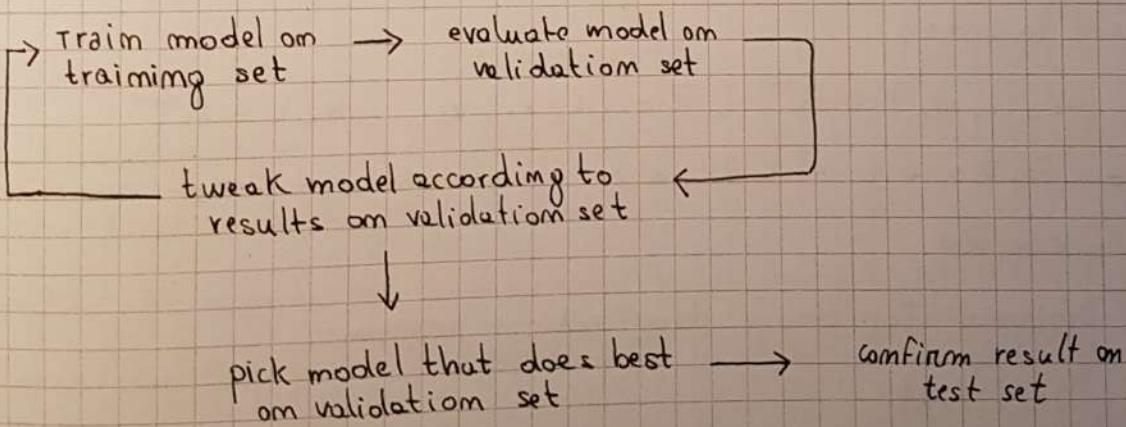
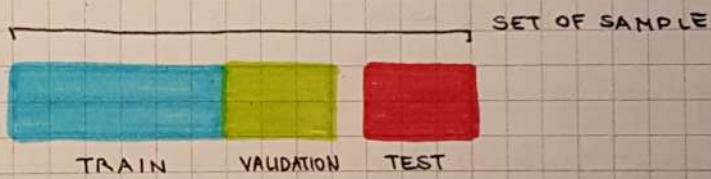
BIAS COMPLEXITY TRADE-OFF

H large \Rightarrow decrease E_{app} , increase E_{est} \Rightarrow overfitting

H small \Rightarrow decrease E_{est} , increase E_{app} \Rightarrow underfitting

TRAIN, TEST AND VALIDATION SETS

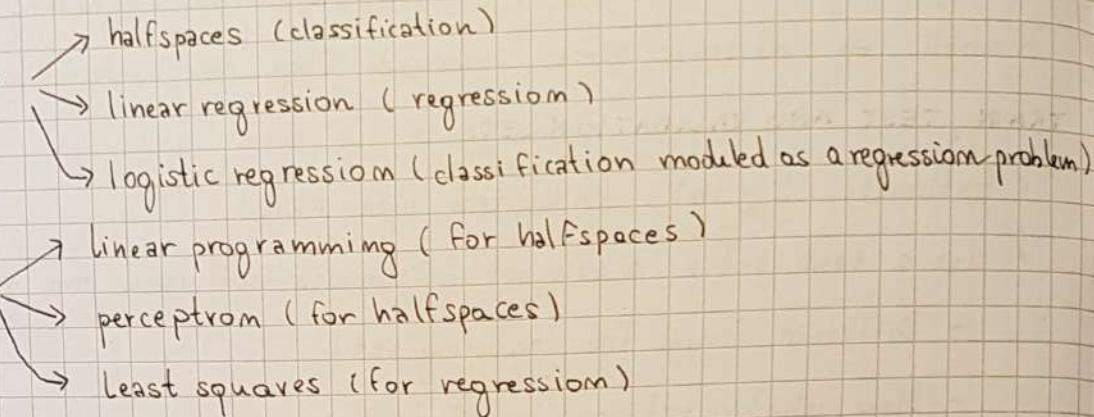
- Estimate $L_D(h)$ (h selected with ERM)
- Use a test set: a new set of samples not used for picking h
 - Different from the training set!
 - More reliable estimation of $L_D(h)$.
 - The test must not be looked at until we have picked our final hypothesis!
 - In practice: 1 set of sample splitted into training set and test set.
- Sometimes: training set divided into training set and validation set.
 - The validation set is used for selecting the hyper-parameters of the algorithm.



LINEAR PREDICTORS (cap.9)

ERM approach.

hypothesis class



Class of AFFINE FUNCTIONS :

$$L_d = \{ h_{\bar{w}, b} : \bar{w} \in \mathbb{R}^d, b \in \mathbb{R} \}$$

$$\text{where } h_{\bar{w}, b}(\bar{x}) = \langle \bar{w}, \bar{x} \rangle + b = \left(\sum_{i=1}^d w_i x_i \right) + b \quad \text{bias}$$

each member of L_d is a function $\bar{x} \rightarrow \langle \bar{w}, \bar{x} \rangle + b$.

\Rightarrow dimensionality of parameters space : $d+1$

hypothesis class : $\phi: \mathbb{R} \rightarrow \mathcal{Y}$ on L_d .

- Binary classification: $\mathcal{Y} = \{-1, 1\} \rightarrow \phi(z) = \text{sigm}(z)$

- Regression: $\mathcal{Y} = \mathbb{R} \rightarrow \phi(z) = z$

More convenient to write :

$$\bar{w}' = (b, w_1, w_2, \dots, w_d) \in \mathbb{R}^{d+1}, \bar{x}' = (1, x_1, x_2, \dots, x_d) \in \mathbb{R}^{d+1}$$

$$\Rightarrow h_{\bar{w}, b}(\bar{x}) = \langle \bar{w}, \bar{x} \rangle + b = \langle \bar{w}', \bar{x}' \rangle \quad \text{now becomes a linear function}$$

$$\text{thus } h_{\bar{w}}(\bar{x}) = \langle \bar{w}, \bar{x} \rangle$$

HALFSPACES Designed for binary classification problems.

$$\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = \{-1, 1\}$$

The class of halfspaces is defined as :

$$HS_d = \text{sigm} \circ L_d = \{ \bar{x} \mapsto \text{sigm}(h_{\bar{w}, b}(\bar{x})) : h_{\bar{w}, b} \in L_d \} \quad \text{hypothesis class}$$

Realizability : linearly separable.

• LINEAR PROGRAMMING

Maximize a linear function subject to linear inequalities.
TARGET:

$$\text{Find } \max_{\bar{w} \in \mathbb{R}^d} \langle \bar{u}, \bar{w} \rangle \quad \text{subject to } A \bar{w} \geq \bar{v}$$

$$\bar{w} \in \mathbb{R}^d \text{ unknowns, } A \in m \times d \text{ matrix, } \bar{v} \in \mathbb{R}^m, \bar{u} \in \mathbb{R}^d$$

Exist efficient linear programming solvers.

Empirical Risk Minimization (ERM) For halfspaces in the realizable case can be expressed as a linear program. In fact:

$$S = \{(\bar{x}_i, y_i)\}_{i=1}^m \text{ training set}$$

- realizable assumption $\Rightarrow L_S(h) = 0 \Rightarrow$ we are looking for \bar{w} such that $y_i \langle \bar{w}, \bar{x}_i \rangle > 0 \quad \forall i = 1, \dots, m$

- There exist also a vector that satisfy $y_i \langle \bar{w}, \bar{x}_i \rangle \geq 1 \quad \forall i = 1, \dots, m$:

proof. w^* satisfy $y_i \langle \bar{w}^*, \bar{x}_i \rangle > 0 \quad \forall i$, $\hat{w} = \frac{w^*}{\gamma}$ with $\gamma = \min_i y_i \langle w^*, \bar{x}_i \rangle$

$$\Rightarrow \forall i \quad y_i \langle \hat{w}, \bar{x}_i \rangle = \frac{1}{\gamma} y_i \langle \bar{w}^*, \bar{x}_i \rangle \geq 1$$

exist a vector such that

$$\Rightarrow y_i \langle \bar{w}, \bar{x}_i \rangle \geq 1 \quad \forall i$$

- Model ERM with linear programming:

$A \in \mathbb{R}^{m \times d}$ matrix such that $A_{ij} = y_i x_{ij} \quad x_{ij} = j\text{-element of the vector } x_i$

$$y_i \langle \bar{w}, \bar{x}_i \rangle \geq 1 \implies A \bar{w} \geq \bar{v} \quad \text{LP constraint } \bar{v} = (1, \dots, 1) \in \mathbb{R}^m$$

$\bar{u} = (0, \dots, 0) \in \mathbb{R}^m$ "dummy target" for $\max_{\bar{w} \in \mathbb{R}^d} \langle \bar{u}, \bar{w} \rangle$

PERCEPTRON

Iterative algorithm.

TARGET: find separating hyperplane.

\rightarrow vector \bar{w} represents separating hyperplane

At each step focus on a missclassified sample and guide the algorithm to be "more correct" on it.

INPUT: training set $(x_1, y_1), \dots, (x_m, y_m)$

initialize $w^{(0)} = (0, \dots, 0)$

For $t = 1, 2, \dots$ do

if $\exists i$ s.t. $y_i \langle w^{(t)}, x_i \rangle \leq 0$ then $w^{(t+1)} = w^{(t)} + y_i x_i$

else return $w^{(t)}$

Note that: $y_i \langle w^{(t+1)}, x_i \rangle = y_i \langle w^{(t)} + y_i x_i, x_i \rangle = y_i \langle w^{(t)}, x_i \rangle + \|x_i\|^2$

Realizability assumption corresponds to linearly separable data.

Theorem

$(x_1, y_1), \dots, (x_m, y_m)$ is linearly separable

$$B = \min \{ \|w\| : y_i \langle w, x_i \rangle \geq 1 \quad \forall i = 1, \dots, m \}$$

$$R = \max_i \|x_i\|$$

The Perceptron algorithm stops after at most $(RB)^2$ iterations, and when it stops it holds that $\forall i \in [m], y_i \langle w^{(t)}, x_i \rangle > 0$.

Proof.

1. Define:

w^* vector achieving the min in B

T number of iterations before stopping

2. Consider:

$$\frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \|w^{(T+1)}\|} \leq 1$$

3. We need to demonstrate: $\frac{\sqrt{T}}{RB} \leq \frac{\langle w^*, w^{(T+1)} \rangle}{\|w^*\| \|w^{(T+1)}\|} \leq 1 \Rightarrow T \leq (RB)^2$

4. Proceed in two parts:

a) Numerator: demonstrate $\langle w^*, w^{(T+1)} \rangle \geq T$

- First iteration: $w^{(1)} = (0, \dots, 0) \rightarrow \langle w^*, w^{(1)} \rangle = 0$

- at each step: $\langle w^*, w^{(t+1)} \rangle - \langle w^*, w^{(t)} \rangle \geq 1$

- after T iterations: $\langle w^*, w^{(T+1)} \rangle \geq T$

$$\text{In fact } \langle w^*, w^{(T+1)} \rangle = \sum_{t=1}^T (\langle w^*, w^{(t+1)} \rangle - \langle w^*, w^{(t)} \rangle) =$$

$$= \sum_{t=1}^T \langle w^*, w^{(t+1)} - w^{(t)} \rangle = \sum_{t=1}^T \langle w^*, y_i x_i \rangle \geq T$$

b) Denominator: demonstrate $\|w^*\| \|w^{(T+1)}\| \leq \sqrt{T} RB$

$$\begin{aligned} \|w^{(t+1)}\|^2 &= \sum_i \|w^{(t)} + y_i x_i\|^2 = \|w^{(t)}\|^2 + 2 \sum_i y_i \langle w^{(t)}, x_i \rangle + \sum_i \|x_i\|^2 \\ &\leq \|w^{(t)}\|^2 + R^2 \end{aligned}$$

$$\Rightarrow \|w^{(T+1)}\| = \sqrt{\sum_{t=1}^T (\|w^{(t+1)}\|^2 - \|w^{(t)}\|^2)} = \sum_{t=1}^T (\|w^{(t)} + y_i x_i\|^2 - \|w^{(t)}\|^2) =$$

$$= \sum_{t=1}^T (2 \sum_i y_i \langle w^{(t)}, x_i \rangle + \|x_i\|^2) \leq TR^2$$

- Rem
- LINEAR
- Estim
- value
- $X \in$
- Final
- impu
- hy po
- Loss
- Empirical

The solution

5) Combining we obtain:

$$\frac{\langle \bar{w}^{(T+1)}, w^* \rangle}{\|\bar{w}^*\| \|\bar{w}^{(T+1)}\|} \geq \frac{T}{B\sqrt{R}} = \frac{\sqrt{T}}{BR}$$

□

Remark

- Convergence is guaranteed!
- Convergence depends on B , which can be exponential in d

LINEAR REGRESSION

Estimate the relation between some explanatory variables and some real valued outcome.

- $X \in \mathbb{R}^d, Y \in \mathbb{R}$
- Find $h \in H_{\text{reg}} : \mathbb{R}^d \rightarrow \mathbb{R}$ that best approximates the relation between input and output, where

hypothesis class $H_{\text{reg}} = L_d = \{\bar{x} \mapsto \langle \bar{w}, \bar{x} \rangle + b : \bar{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$

- Loss Function: L2 commonly used

$$l(h, (\bar{x}, y)) = (h(\bar{x}) - y)^2$$

- Empirical risk Function: $L_s(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2$ MEAN SQUARED ERROR

LEAST SQUARE

It's the algorithm that solves the ERM problem for linear regression predictions with the squared loss.

ERM problem is to find

$$\arg \min_{\bar{w}} L_s(h_{\bar{w}}) = \arg \min_{\bar{w}} \frac{1}{m} \sum_{i=1}^m (\langle \bar{w}, \bar{x}_i \rangle - y_i)^2$$

Solve the problem: calculate gradient w.r.t. \bar{w} and set to \emptyset :

$$\frac{\partial L}{\partial \bar{w}} = \frac{2}{m} \sum_{i=1}^m (\langle \bar{w}, \bar{x}_i \rangle - y_i) \bar{x}_i = 0$$

$$\Rightarrow A \bar{w} = \bar{b} \quad \text{with} \quad A = \left(\sum_{i=1}^m x_i x_i^\top \right), \quad \bar{b} = \sum_{i=1}^m y_i \bar{x}_i$$

The solution is: $\bar{w} = A^{-1} \bar{b}$

POLYNOMIAL REGRESSION

- Find the one dimensional polynomial of degree m that better predicts the data

$$p(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_m x^m$$

$$(a_0, \dots, a_m) \in \mathbb{R}^{m+1}$$

$$\mathcal{H}_{\text{poly}}^m = \{x \mapsto p(x)\} \text{ hypothesis class}$$

N.B. $X = \mathbb{R}$, $Y = \mathbb{R} \Rightarrow$ regression problem

\Rightarrow reduce the problem to a m -dimensional linear regression using the map

$$\psi: \mathbb{R} \rightarrow \mathbb{R}^{m+1} \quad \text{s.t. } \psi(x) = (1, x, x^2, \dots, x^m)$$

We obtain

$$\langle \vec{a}, \psi(x) \rangle = a_0 + a_1 x + a_2 x^2 + \dots + a_m x^m$$

\Rightarrow Find the optimal vector of coefficients \vec{a} by using the Least Squares algorithm.

LOGISTIC REGRESSION

Reframe a classification problem as a regression one.

TARGET: learn a function $h: \mathbb{R}^d \rightarrow [0, 1]$

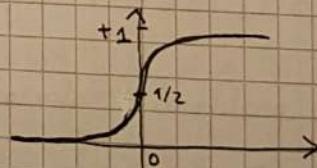
It's used for classification tasks: we can interpret $h(\bar{x})$ as the probability that the label of \bar{x} is 1.

hypothesis class: $H: \Phi_{\text{sig}} \circ L_d = \{\bar{x} \mapsto \Phi_{\text{sig}}(\langle \bar{w}, \bar{x} \rangle) : \bar{w} \in \mathbb{R}^d\}$

where

$\Phi_{\text{sig}}: \mathbb{R} \rightarrow [0, 1]$ sigmoid function over linear function class

$$\Phi_{\text{sig}}(z) = \frac{1}{1 + e^{-z}}$$



Loss function: $\ell(h_{\bar{w}}, (\bar{x}, y)) = \log(1 + \exp(-y \langle \bar{w}, \bar{x} \rangle))$

ERM problem: $\arg \min_{\bar{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle \bar{w}, \bar{x}_i \rangle))$

MLE

Maximum Likelihood estimation (MLE) is a statistical approach for finding the parameters that maximize the joint probability of a given dataset assuming a specific parametric probability function.

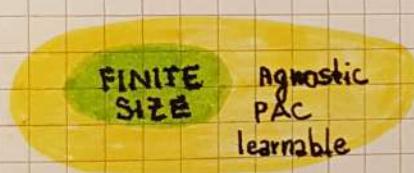
\Rightarrow is equivalent to ERN solution for logistic regression

- $S = ((\bar{x}_1, y_1), \dots, (\bar{x}_m, y_m))$ with each (\bar{x}_i, y_i) i.i.d. from some probability distribution.
- $P[S(\theta)]$ (likelihood of data given parameters)
- log-likelihood : $L(S; \theta) = \log(P[S(\theta)])$
- maximum likelihood estimator : $\hat{\theta} = \arg \max_{\theta} L(S; \theta)$

THE VC- DIMENSION (cap.6)

AIM: answer at the question: which hypothesis class are PAC learnable?

Simplification: Focus on binary classification and 0-1 loss.



if $|H| < \infty \Rightarrow H$ is PAC learnable.

if $|H| = \infty$?

\Rightarrow we'll demonstrate that the finite size is a sufficient but not necessary condition

Definition (Restriction of H to C)

let H be a class of functions from X to $\{0,1\}$ and let $C = \{c_1, \dots, c_m\} \subset X$.

The restriction H_C of H to C is :

$$H_C = \{[h(c_1), \dots, h(c_m)] : h \in H\}$$

namely, the set of functions from C to $\{0,1\}$ that can be derived from H . We represent each function from C to $\{0,1\}$, as a vector in $\{0,1\}^{|C|}$.

Def (Shattering)

A hypothesis class shatters a finite set $C \subset X$ if the restriction of H to C is the set of all functions from C to $\{0,1\}$. That is, $|H_C| = 2^{|C|}$.

• AXIS A

Corollary (of NO-FREE LUNCH)

Let H be a hypothesis class of functions from X to $\{0,1\}$. Let m be a training size. Assume that there exists a set $C \subset X$ of size $2m$ that is shattered by H . Then, for any learning algorithm, A , there exist a distribution D over $X \times \{0,1\}$ and a predictor $h \in H$ such that $L_D(h) = 0$ but with probability of at least $1/7$ over the choice of $S \sim D^m$ we have that $L_{D_S}(A(S)) \geq 1/8$.

Remark The corollary tell us that if H shatters some set C of size $2m$, then we cannot learn H using m samples.

Definition (VC-dimension)

The VC-dimension of a hypothesis class H , denoted $\text{VCdim}(H)$, is the maximal size of a set $C \subset X$ that can be shattered by H .

NOTE: if H can shatter sets of arbitrarily large size, then we say that

$$\text{VCdim}(H) = +\infty.$$

Theorem

Let H be a class of infinite VC-dimension. Then, H is not PAC learnable.

Proof. consequence of previous NFL corollary, there will always exist a subset of size $2m$ that is shattered for any m .

Remark In the case of finite class hypothesis: $|H| < +\infty \Rightarrow \text{VCdim}(H) \leq \log_2 |H|$.

EXAMPLES

To show that $\text{VCdimension}(H) = d$ we need to show that:

1. $\text{VCdim}(H) \geq d$: there exists a set C of size d which is shattered by H .
2. $\text{VCdim}(H) < d+1$: every set of size $d+1$ is not shattered by H .

• THRESHOLD FUNCTIONS

$$H = \{h_a : a \in \mathbb{R}\} \quad \text{with } h_a : \mathbb{R} \rightarrow \{0,1\} \quad h_a(x) = \mathbb{1}[x < a] = \begin{cases} 1 & \text{if } x < a \\ 0 & \text{if } x \geq a \end{cases}$$

$$\text{VCdim}(H) = 1$$

• INTERVALS

$$H = \{h_{a,b} : a, b \in \mathbb{R}, a < b\} \quad \text{with } h_{a,b} : \mathbb{R} \rightarrow \{0,1\} \quad h_{a,b}(x) = \mathbb{1}[x \in (a, b)] = \begin{cases} 1 & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

$$\text{VCdim}(H) = 2$$

• AXIS A

$H = \{f\}$

with

VCdim

Remark

↓

• HALF

• VC di

• VC di

FUNDAMENTAL

Theorem

Let H

the loss

1. N

2. A

5.

6. H

Notes

1. We

2. 3

3. 2

4. 4

5. T

The p

• I

is

• AXIS ALIGNED RECTANGLES

$$H = \{h_{(a_1, a_2; b_1, b_2)} : a_1, a_2, b_1, b_2 \in \mathbb{R}, a_1 \leq a_2, b_1 \leq b_2\}$$

with $h_{(a_1, a_2, b_1, b_2)}(x_1, x_2) = \begin{cases} 1 & \text{if } a_1 \leq x_1 \leq a_2, b_1 \leq x_2 \leq b_2 \\ 0 & \text{otherwise} \end{cases}$

$$VC(\dim(H)) = 4$$

Remark. In general For K-dimension space : $VC(\dim(H)) = K+1$

• HALFSPACE HYPOTHESIS CLASS

- $VC(\dim)$ of homogeneous halfspaces in \mathbb{R}^d is d .
- $VC(\dim)$ of non-homogeneous halfspaces in \mathbb{R}^d is $d+1$.

FUNDAMENTAL THEOREM OF PAC LEARNING

Theorem (The fundamental theorem of statistical learning)

Let H be a hypothesis class of functions from a domain X to $\{0, 1\}$ and let the loss function be the 0-1 loss. Then, the following are equivalent:

1. H has the uniform convergence property.
2. Any ERM rule is a successful agnostic PAC learner for H .
3. H is agnostic PAC learnable.
4. H is PAC learnable.
5. Any ERM rule is successful PAC learner for H .
6. H has a finite VC-dimension.

Notes on the demonstration:

1. We've already seen that $1 \Rightarrow 2 \Rightarrow 3$.
2. $3 \Rightarrow 4$ is trivial
3. $2 \Rightarrow 5$ also trivial
4. $4 \Rightarrow 6$ and $5 \Rightarrow 6$ (infinite VC is not PAC learnable)
5. The key point is how to close the loop $6 \Rightarrow 1$!

The proof of $6 \Rightarrow 1$ can be divided in two main parts:

- If $VC(\dim(H)) = d$, then even though $|H|$ might be infinite, when restricting $|H|$ to a finite set C , it's 'effective size' $|H_C|$ is only $O(|C|^d)$. That is, $|H_C|$ grows polynomially rather than exponentially with $|C|$. — SAVERS LEMMA

Recall that infinite hypothesis classes enjoy the uniform convergence property. This result can be generalized by showing that uniform convergence holds whenever the hypothesis class has a "small effective size" (i.e., classes for which $|H|$ grows polynomially with $|C|$).

Theorem (The fundamental theorem of statistical learning - quantitative)

Let H be a hypothesis class of functions from a domain X to $\{0,1\}$ and let the loss function be the 0-1 loss. Assume that $\text{VC dim}(H) = d < \infty$. Then, there are absolute constants C_1, C_2 such that:

1. H has the uniform convergence property with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_H^{vc}(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

2. H is agnostic PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon^2} \leq m_H(\epsilon, \delta) \leq C_2 \frac{d + \log(1/\delta)}{\epsilon^2}$$

3. H is PAC learnable with sample complexity

$$C_1 \frac{d + \log(1/\delta)}{\epsilon} \leq m_H(\epsilon, \delta) \leq C_2 \frac{d \log(1/\epsilon) + \log(1/\delta)}{\epsilon}$$

Definition (Growth function)

Let H be a hypothesis class. Then the growth function of H , denoted $\gamma_H : \mathbb{N} \rightarrow \mathbb{N}$, is defined as

$$\gamma_H(m) = \max_{C \subset X : |C|=m} |H_C|$$

In words, $\gamma_H(m)$ is the number of different functions from a set C of size m to $\{0,1\}$ that can be obtained by restricting H to C .

Notes:

- $\gamma_H(m)$ is monotonically nondecreasing.
- $\gamma_H(m)$ is distribution free (does not depend on D)
- What is the upper bound to $\gamma_H(m)$?

Remark. If $\text{VC dim}(H) = d$, then for any $m \leq d$ we have $\gamma_H(m) = 2^m$, therefore H induces all possible functions from C to $\{0,1\}$.

Lemma (Sauer-Shelah-Perles)

Let H be a hypothesis class with $\text{VCdim}(H) \leq d < \infty$. Then, for all m ,

$$\gamma_H(m) \leq \sum_{i=0}^d \binom{m}{i}.$$

$$\text{N.B. } \binom{m}{i} = \frac{m!}{i!(m-i)!}$$

In particular, if $m > d+1$, then

$$\gamma_H(m) \leq \left(\frac{em}{d}\right)^d$$

N.B. In this case, it grows polynomially with m !

Remark. For $m > d+1$ the bound is much smaller than 2^m !

Theorem

Let H be a class and let γ_H be its growth function. Then, for every D and every $\delta \in (0,1)$, with probability of at least $1-\delta$ over the choice of $S \sim D^m$ we have

$$|L_D(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\gamma_H(2m))}}{8\sqrt{2m}}$$

Proof idea. If we take two independent samples S, S' of size m call them, then $|L_S(h) - L_{S'}(h)|$ should look like an estimate for a hypothesis class of size $\gamma_H(2m)$.

A FEW NOTES ON NON-UNIFORM LEARNABILITY (cap. 7) no exam

Def

A hypothesis class H is mom-uniformly learnable if there exist a learning algorithm, A , and a function $m_H^{\text{NUL}} : (0,1)^2 \times H \rightarrow \mathbb{N}$ such that, for every $\epsilon, \delta \in (0,1)$ and for every $h \in H$, if $m \geq m_H^{\text{NUL}}(\epsilon, \delta, h)$ then for every distribution D , with probability of at least $1-\delta$ over the choice of $S \sim D^m$, it holds that

$$L_D(A(S)) \leq L_D(h) + \epsilon$$

Key difference: m depends on h !

Recall the definition of AGNOSTIC PAC learnability.

\rightarrow we note that this definition implies that $\forall h \in H : L_D(A(S)) \leq L_D(h) + \epsilon$

\Rightarrow AGNOSTIC PAC learnability \Rightarrow mom-uniformly learnable

Theorem

A hypothesis class H of binary classifiers is mom uniformly learnable if and only if it is a countable union of agnostic PAC learnable hypothesis classes ($H = \bigcup_{m \in \mathbb{N}} H_m$, each H_m agnostic PAC).

Theorem

Let H be a hypothesis class that can be written as a countable union of hypothesis classes $H = \bigcup_{m \in \mathbb{N}} H_m$, where each H_m enjoys the uniform convergence property. Then H is mom uniformly learnable.

Non-uniform learning can be achieved through the STRUCTURAL RISK MINIMIZATION (SRM) paradigm where we jointly minimize the empirical risk and a cost function depending on the complexity of the selected hypothesis

$$h \in \operatorname{argmin}_{h \in H} [L_S(h) + \epsilon_{\text{mdl}}(m, w(m), \cdot, \delta)]$$

Assume a description language for the hypothesis class H . Denote with $|h|$ the length of description of hypothesis $h \in H$.

Theorem (Minimum description length) (MDL)

Let H be a hypothesis class, let $d: H \rightarrow \{0, 1\}$ be a prefix-free description language for H . Then:

- For every sample size m
- For every $\delta > 0$
- For every probability distribution D

with probability $\geq 1 - \delta$ (over the choice $S \sim D^m$) we have

$$\forall h \in H = L_D(h) \leq L_S(h) + \sqrt{\frac{|h| + \log(2/\delta)}{2m}}$$

The minimum description length (MDL) represents an alternative learning paradigm for H (in place of ERM):

prior Knowledge:

- H is countable hypothesis class
- H is described by a prefix-free language over $\{0,1\}$
- for every $h \in H$: $|h| =$ length of representation of h

input:

training set $S \sim D^m$, confidence δ

output:

$$\text{MDL}(S) : h \in \underset{h \in H}{\operatorname{argmin}} \left[L_S(h) + \sqrt{\frac{|h| + \delta m (2/\delta)}{2m}} \right]$$

"bias" complexity

Occam's Razor:

"A short explanation (that is, a hypothesis that has a short length) tends to be more valid than a long explanation."

MODEL SELECTION AND VALIDATION (cap. 11)

Aim: choose the best algorithm for the particular problem.

There are different algorithms. Algorithms have parameters / design choices.

How to select the best?

Two approaches:

- STRUCTURAL RISK MINIMIZATION (SRM)
 - limited practical applicability
 - not part of the course
- Use a validation test

MODEL SELECTION WITH SRM (no exam)

- $H = \bigcup_{m \in \mathbb{N}} H_m$
- H_m : uniform convergence property.
- sample complexity of $m_{H_m}^{VC}(\epsilon, \delta) \leq \frac{g(m) \log(1/\delta)}{\epsilon^2}$

where $g: \mathbb{N} \rightarrow \mathbb{R}$ monotonically increasing

Bound: with probability $\geq 1 - \delta$, for every $d \in \mathbb{N}$ and $h \in H_d$

$$L_D(h) \leq L_S(h) + \sqrt{\frac{g(d)(\log(1/\delta)) + 2 \log d + \log(\pi^2/\delta)}{m}}$$

SRM rule: pick d and $h \in H_d$ minimizing

$$L_s(h) + \sqrt{\frac{g(d)(\log(1/\delta) + 2\log(d) + \log(\pi^2/6))}{m}}$$

In many practical cases the upper bound is pessimistic.

VALIDATION SET

Idea: divide the training set in 2 parts, use the first to pick an hypothesis, and the second (not used to train) to estimate its true error.

Assume we have picked a predictor h (e.g. by ERM rule on H_d):

- V : set of m_V samples from \mathcal{D} not used for training (validation set)
- L_V : loss computed on V (loss in $[0,1]$)

Theorem (1)

Let h be some predictor and assume that the loss function is in $[0,1]$.

Then, for every $\delta \in (0,1)$, with probability of at least $1-\delta$ over the choice of a validation set V of size m_V , we have

$$|L_V(h) - L_D(h)| \leq \sqrt{\frac{\log(2/\delta)}{2m_V}}$$

Proof. using Hoeffding's inequality.

Remark. The bound in the last theorem does not depend on the algorithm or the training set used to construct h .

The bound based on the validation set is more accurate:

- there is no constant term
- does not depend on VC dimension (because validation samples have not been used for training).
- choose final hypothesis by ERM over the validation set

To illustrate this points:

- suppose that h was obtained by applying an ERM predictor with respect to a hypothesis class of VC-dimension d , over a training set of m samples. Then, from the fundamental theorem of learning we obtain the bound

$$L_D(h) \leq L_s(h) + \sqrt{C \frac{d + \log(1/\delta)}{m}}$$

• In contrast, from the last theorem we obtain the bound

$$L_D(h) \leq L_V(h) + \sqrt{\frac{\log(2/\delta)}{2m_V}}$$

What if we have one or more parameters with values in \mathbb{R} ?

1. Start with a rough grid of values.
2. Plot the corresponding model-selection curve.
3. Based on the curve, zoom in to the correct region.
4. Restart from (1) with a finer grid.

Remark. The empirical risk on the validation set is not an estimate of the true risk, in particular if we choose among many models!

Question: how can we estimate the true risk after model selection?

Summary

We have to choose among multiple possible hypothesis set H_i .

Approach: split the data in 3 parts.

1. TRAINING SET: used to learn the best model h_i inside each class H_i .
2. VALIDATION SET: used to pick one hypothesis h^* from h_1, h_2, \dots .
3. TEST SET: used to estimate the true risk $L_D(h^*)$.

Note that:

- the estimate from the test set has the guarantees provided by the proposition on estimate of $L_D(h^*)$ for one class.
- The test set is not involved in the choice of h^* .
- If after using the test set to estimate $L_D(h^*)$ we decide to choose another hypothesis (because we have seen $L_D(h^*)$) we cannot use the test set again to estimate $L_D(h^*)$!



K-FOLD CROSS VALIDATION

When data is not plentiful, we cannot afford to drop part of it to build the validation set. We use the K-fold cross validation.

K-Fold cross validation =

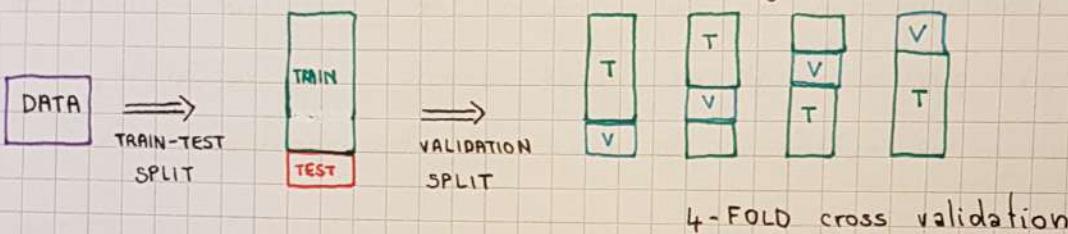
1. Partition (training) set of m samples into K folds of size m/K .

2. For each fold:

- train on union of the other folds

- estimate error (for learned hypothesis) on the selected folds.

3. Estimate of the true error as the average of the estimated errors.



Remark If $K = m \Rightarrow$ leave-one-out (LOO) case.

K-fold cross validation is often used for model selection, and once the best parameter is chosen, the algorithm is retrained using this parameter on the entire training set.

Pseudocode :

input :

training set $S = (\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)$

set of parameter values Θ

learning algorithm A

integer K

partition: S into S_1, S_2, \dots, S_K

for each $\theta \in \Theta$

for $i = 1, \dots, K$

$$h_{i,\theta} = A(S \setminus S_i; \theta)$$

$$\text{error}(\theta) = \frac{1}{K} \sum_{i=1}^K L_{S_i}(h_{i,\theta})$$

output :

$$\theta^* = \arg \min_{\theta} [\text{error}(\theta)]$$

$$h^* = A(S; \theta^*)$$

WHAT TO DO IF LEARNING FAILS

Case: good training/validation errors, but results on test set are bad!

Main approaches to solve the problem:

- Get a larger sample.
- Change the hypothesis class by:
 - enlarging it
 - reducing it
 - completely changing it
 - changing the parameters you consider
- Change the feature representation of the data.
- Change the optimization algorithm used to apply your learning rule.

In order to find the best remedy, it is essential first to understand the cause of the bad performance.

Recall:

h_s = learned predictor

- $L_D(h^*)$ APPROXIMATION ERROR (true error of best hypothesis in H)
- $L_D(h_s) - L_D(h^*)$ ESTIMATION ERROR
 $\rightarrow L_D(h_s)$
 (difference between the true error of best hypothesis in H and true error of ERM solution)
- $L_S(h_s)$ TRAINING ERROR (empirical error of ERM solution on training set S)
- $L_V(h_s)$ VALIDATION ERROR (error on validation set V of ERM solution)

We decomposed the true error of the learned predictor into:

$$L_D(h_s) = L_D(h^*) + (L_D(h_s) - L_D(h^*))$$

using train and validation errors

$$L_D(h_s) = \underbrace{(L_D(h_s) - L_V(h_s))}_{\text{can be bounded by theorem (1)}} + \underbrace{(L_V(h_s) - L_S(h_s))}_{\text{if is large } \rightarrow \text{overfitting}} + L_S(h_s)$$

Remark

- The approximation error of the class does not depend on the sample size or on the algorithm being used. It only depends on the distribution D and on the hypothesis class H .
- The estimation error of the class does depend on the sample size.

The following steps should be applied:

1. If learning involves parameter tuning, plot the model-selection curve to make sure that you tuned the parameters appropriately.
2. If the training error is excessively large consider enlarging the hypothesis class, completely change it, or change the feature representation of the data.
- ⇒ 3. If the training error is small, plot learning curves and try to deduce from them whether the problem is estimation error or approximation error.
4. If the approximation error seems to be small enough, try to obtain more data. If this is not possible, consider reducing the complexity of the hypothesis class.
5. If the approximation error seems to be large as well, try to change the hypothesis class or the feature representation of the data completely.

$$L_s(h_s) = (L_s(h_s) - L_s(h^*)) + (L_s(h^*) - L_p(h^*)) + L_p(h^*)$$

Remark For (3): if the validation error seems to decrease (the error is large but the two curves get closer):

- get more data (if possible)
- otherwise reduce complexity of H

If the validation error remains large:

- change H
- change feature representation of the data

REGULARIZATION AND STABILITY (cap. 13)

The new learning paradigm we introduce in this chapter is called

REGULARIZED LOSS MINIMIZATION (RLM).

Key idea: jointly minimize empirical risk and a regularization function.

- hypothesis h : defined by a vector $\vec{w} = (w_1, \dots, w_d)^T \in \mathbb{R}^d$
(e.g. coefficients of a linear model, weights in a neural network, etc.)
- Regularization Function: $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}$, function of \vec{w} .

• Regularized loss minimization (RLM) = select \vec{w} from

$$\operatorname{argmin}_{\vec{w}} (L_s(\vec{w}) + R(\vec{w}))$$

• $L_s(\vec{w})$: standard loss as previously

• $R(\vec{w})$: regularization term (measures in some way the "complexity" of the found solution)

The regularization term balances between low empirical risk and aiming at less complex hypothesis.

It is possible to view the extra term as a "stabilizer".

TIKHONOV REGULARIZATION

Define a function R using the ℓ_2 norm:

$$R(\vec{w}) = \lambda \|\vec{w}\|^2 = \lambda \sum_{i=1}^d w_i^2 \quad \text{where } \lambda > 0 \text{ is a scalar}$$

The output of the function R is a real positive number.

This yields the learning rule:

$$A(S) = \operatorname{argmin}_{\vec{w}} (L_s(\vec{w}) + \lambda \|\vec{w}\|^2)$$

Remark $\|\vec{w}\|^2$ measures the "complexity" of the hypothesis defined by \vec{w} .

Instead, λ controls the trade-off between low empirical risk or high complexity (risk of overfitting).

RIDGE REGRESSION

Linear regression with squared loss + Tikhonov regularization

• Linear regression with squared loss:

Find \vec{w} that minimizes the squared loss

$$\vec{w} = \operatorname{argmin}_{\vec{w}} \sum_{i=1}^m (\langle \vec{w}, \vec{x}_i \rangle - y_i)^2$$

• Ridge regression: find \vec{w} that minimizes

$$\vec{w} = \operatorname{argmin}_{\vec{w}} \left(\lambda \|\vec{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\langle \vec{w}, \vec{x}_i \rangle - y_i)^2 \right)$$

Closed form solution:

aim: Find optimal \vec{w}

1. Compute gradient w.r.t. \vec{w} and set to 0

$$2\lambda \|\vec{w}\|^2 + \frac{1}{m} \sum_{i=1}^m (\langle \vec{w}, \vec{x}_i \rangle - y_i) \vec{x}_i = 0$$

Proposition

If al

E_s

Proof.

that

E_s

On t

E_s

The

The p

only i

→ Stable

TIKHONOV

Def (

Lemmas

1. Th

2. I

3. I

Def (

Let

V → W₁

Remark

More or

derivative

2. Define

$$A = \left(\sum_{i=1}^m \vec{x}_i \vec{x}_i^\top \right)$$

$$\vec{b} = \sum_{i=1}^m y_i \vec{x}_i$$

3. The solution is

$$\vec{w} = (2\lambda m \mathbb{1} + A)^{-1} \vec{b}$$

Remark Tikhonov regularization makes the learner stable w.r.t. small perturbations of the training set. This turns leads to small bounds on generalization error.

STABILITY

Informally: an algorithm A is stable if a small change of the training data (i.e., its input) S will lead to a small change of its output hypothesis.

Let A be a learning algorithm, let S = (z₁, ..., z_m) be a training set of m samples and let A(S) denote the output of A.

Let S⁽ⁱ⁾ be the training set obtained by replacing the i'th example of S with z': namely, S⁽ⁱ⁾ = (z₁, ..., z_{i-1}, z', z_{i+1}, ..., z_m).

In our informal definition =

- "a small change of the input" means that we feed A with S⁽ⁱ⁾ instead of with S. That is, we only replace one training sample.
- "a small change of its output hypothesis" means a small change in the loss → ON-AVERAGE-REPLACE-ONE-STABLE (OAROS) algorithms.

Def (ON-AVERAGE-REPLACE-ONE-STABLE)

Let E: N → ℝ be a monotonically decreasing function. We say that a learning algorithm A is on-average-replace-one-stable (OAROS) with rate E(m) if for every distribution D:

$$\mathbb{E}_{(S, z') \sim D^{m+1}, i \sim U(m)} [\ell(A(S^{(i)}), z_i) - \ell(A(S), z_i)] \leq E(m)$$

Remark U(m) is a uniform distribution over [m]

Proposition

If algorithm A is OAROS with rate $\epsilon(m)$ then:

$$\mathbb{E}_{S \sim D^m} [L_D(A(S)) - L_S(A(S))] \leq \epsilon(m)$$

$L_D(A(S))$ = true risk of output
 $L_S(A(S))$ = empirical risk of output

Proof. Since S and z' are both drawn i.i.d. from D , we have that for every i :

$$\mathbb{E}_S [L_D(A(S))] = \mathbb{E}_{S, z'} [\rho(A(S), z')] = \mathbb{E}_{S, z'} [\rho(A(S^{(i)}), z_i)]$$

On the other hand we have:

$$\mathbb{E}_S [L_S(A(S))] = \mathbb{E}_{S, i} [\rho(A(S), z_i)] \quad (?)$$

The proof follows from the definition of stability. \square

The proposition tells us that a learning algorithm does not overfit if and only if it is om-average-replace-one-stable —

→ Stable rules do not overfit.

TIKHONOV REGULARIZATION AS A STABILIZER

Def (strongly convex functions)

A function f is λ -strongly convex if for all \vec{w}, \vec{u} , and $\alpha \in (0, 1)$ we have

$$f(\alpha \vec{w} + (1-\alpha) \vec{u}) \leq \underbrace{\alpha f(\vec{w}) + (1-\alpha) f(\vec{u})}_{\text{convex}} - \underbrace{\frac{\lambda}{2} \alpha(1-\alpha) \|\vec{w} - \vec{u}\|^2}_{\text{margin ("strongly")}}$$

Lemma

1. The function $f(\vec{w}) = \lambda \|\vec{w}\|^2$ is 2λ -strongly convex.
2. If f is λ -strongly convex and g is convex, then $f+g$ is λ -strongly convex.
3. If f is λ -strongly convex and \vec{u} is a minimizer of f , then, for any \vec{w} ,

$$f(\vec{w}) - f(\vec{u}) \geq \frac{\lambda}{2} \|\vec{w} - \vec{u}\|^2$$

Def (LIPSCHITZNESS)

Let $C \subset \mathbb{R}^d$, a function $f : \mathbb{R}^d \rightarrow \mathbb{R}^K$ is p -Lipschitz over C if

$$\forall \vec{w}_1, \vec{w}_2 \in C \text{ we have that } \|f(\vec{w}_1) - f(\vec{w}_2)\| \leq p \|\vec{w}_1 - \vec{w}_2\|$$

Remark Intuitively it means that the function cannot change too fast. Moreover, for differentiable functions corresponds to bound on derivative: if derivative bounded by p at any point, the function is p -Lipschitz.

Corollary

Assume that the loss function is convex and p -Lipshitz continuous. Then, the RLM rule with regularizer $\lambda \|\vec{w}\|^2$ is OAROS with rate $\frac{2p^2}{\lambda m}$. It

follows that for the RLM rule:

$$E_{S \sim D^m} [L_D(A(S)) - L_S(A(S))] \leq \frac{2p^2}{\lambda m}$$

CONTROLLING THE FITTING-STABILITY TRADE-OFF

We can rewrite the expected risk of a learning algorithm as

$$E_S [L_D(A(S))] = E_S [L_S(A(S))] + E_S [L_D(A(S)) - L_S(A(S))]$$

- $E_S [L_S(A(S))]$: reflects how well A fits the training set S .
- $E_S [L_D(A(S)) - L_S(A(S))]$: reflects the difference between the true and empirical risks of $A(S)$. It is equivalent to the stability of A .

In Tikhonov regularization, λ controls tradeoff between the two terms

- How do $L_S(A(S))$ and $\|\vec{w}\|^2$ vary as a function of λ ?
→ Larger λ leads to higher empirical risk $L_S(A(S))$.
- How may $E_S [L_D(A(S)) - L_S(A(S))]$ change as a function of λ ?
→ On the other side increasing λ , the stability term decreases.
- How to set λ ?
→ theoretical bound in the book.
→ In practice, validation error is used.

Corollary

Assume that the loss function is convex and p -Lipshitz. Then the RLM rule with the regularization function $\lambda \|\vec{w}\|^2$ satisfies

$$\forall \vec{w}^*, E_S [L_D(A(S))] \leq L_D(\vec{w}^*) + \lambda \|\vec{w}^*\|^2 + \frac{2p^2}{\lambda m}$$

This bound is often called "oracle inequality".

→ If w^* is an hypothesis with low risk, the bound tell us how many examples we need to get close to w^* .

→ But there is a problem...

• need
• it is

STOCHAS

Aim : mee
function

Recall : th
 $\nabla f(\vec{w})$

Idea: th
close to v

• Move im
• Gradient

→ I-

• The th
subgra

GRADIENT
General

GD ALG

Remark

Hypothesi

- $f(\vec{w})$ is
- $\vec{w}^* \in a$

We have

- need to know $\|W^*\|$: in practice not known.
- it is more practical to use a validation set!

STOCHASTIC GRADIENT DESCENT (cap. 14)

Aim: need a general approach to minimize a differentiable convex function $f(\vec{w})$.

Recall: the gradient $\nabla f(\vec{w})$ of a differentiable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is

$$\nabla f(\vec{w}) = \left(\frac{\partial f(\vec{w})}{\partial w_1}, \dots, \frac{\partial f(\vec{w})}{\partial w_d} \right)$$

Idea: the gradient points in the direction of the largest increase of f close to \vec{w} .

- Move in the opposite direction until you find a minima.
- Gradient correspond to first order Taylor approximations.
→ It is good for small steps → need to move step by step.
- The theory can be extended to non-differentiable functions using subgradients.

GRADIENT DESCENT (GD)

General approach to minimize a differentiable convex function $f(\vec{w})$.

GD ALGORITHM:

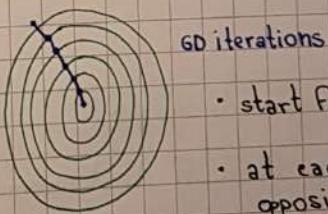
$$\vec{w}^{(0)} \leftarrow \emptyset$$

For $t \leftarrow 0$ to $T-1$ do

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} - \eta \nabla f(\vec{w}^{(t)})$$

return $\vec{w} = \vec{w}^{(T)}$

Remark. $\eta > 0$.



GD iterations

- start from an initial point
- at each step move in direction opposite to the gradient
- stop when solution does not improve or max iterations reached
- get the final point or the one corresponding to minimum value of the objective functions.

Hypothesis:

- $f(\vec{w})$ is a convex p -Lipschitz function
- $\vec{w}^* \in \operatorname{argmin}_{\{\vec{w}: \|\vec{w}\| \leq B\}} f(\vec{w})$

We have the following Corollary:

$\nabla \mathcal{L}(\vec{w})$

Remark. For SGD we have less stable trajectory:

- More "noisy" but could ~~not~~ jump out of local minima.
- Advanced approaches to stabilize, e.g., momentum.
- Sometimes the final point is computed as average of a set of samples to account for fluctuations.
- Better to average only a set of final iterations, not all.
- Alternative approach to get a stable result: use an adaptive step size.

GRADIENT DESCENT: VARIANTS

1. BATCH GRADIENT DESCENT (standard GD): compute the gradient over the complete training set.
2. MINI-BATCH GRADIENT DESCENT: compute the gradient over a small set of K samples.
 - K : parameter, mini-batch size.
 - Trade-off between the two "extreme" cases GD and SGD.
 - Used to train deep neural networks.
3. STOCHASTIC GRADIENT DESCENT (SGD): use a single sample to estimate the gradient.

SGD: APPLICATIONS

Use SGD to solve ML problems:

1. RISK MINIMIZATION (ERM).
2. REGULARIZED LOSS MINIMIZATION (RLM) in details not part of the course.
3. Support vector machines (SVM).
4. Neural networks.

SGD FOR RISK MINIMIZATION

- SGD allows us to minimize $L(\vec{w})$ directly
- All we need is to find an unbiased estimate of the gradient of $L(\vec{w})$, that is a random vector whose conditional expected value is $\mathbb{E}[\nabla L(\vec{w}^{(t)})]$
- Sample a single fresh sample and estimate the gradient with it.

- can be applied to RLM solving its target.

Algorithm: SGD for minimizing $L_D(\vec{w})$

parameters: scalar $\eta > 0$, integer $T > 0$

initialize : $\vec{w}^{(0)} = \vec{0}$

For $t = 1, 2, \dots, T$

sample $z \sim D$

pick $\vec{v}_t \in \nabla L(\vec{w}^{(t)}, z)$

update $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \vec{v}_t$

output $\vec{w} = \frac{1}{T} \sum_{t=1}^T \vec{w}^{(t)}$

We shall now use our analysis of SGD to obtain a sample complexity analysis for learning convex-Lipschitz-bounded problems!

Corollary

Consider a convex-Lipschitz-bounded learning problem with parameters p, B . Then, for every $\epsilon > 0$, if we run the SGD method for minimizing $L_D(\vec{w})$ with a number of iterations (i.e. number of samples)

$$T \geq \frac{B^2 p^2}{\epsilon^2}$$

and with $\eta = \sqrt{\frac{B^2}{p^2 T}}$, then the output of SGD satisfies

$$\mathbb{E}[L_D(\vec{w})] \leq \min_{\vec{w} \in \mathcal{W}} L_D(\vec{w}) + \epsilon$$

SGD: ISSUES

• The selection of the learning rate η is a critical point

→ if it is too small the convergence can be very slow;

→ if it is too large the optimization can be very unstable.

• Simple solution: use adaptive learning rates, e.g.,

↪ progressively reducing the learning rate according to a pre-defined schedule

↪ reduce it when the change in the loss becomes too small

↪ however, these approaches require rules and thresholds to be defined in advance and thus are difficult to adapt to different problems.

- Additionally, the same learning rate applies to all parameters updates
 - the various parameters have different behaviours and the learning rate could be too fast for some and too slow for others.
 - sometimes better not to update all parameters to the same extent, but perform a larger update for some and smaller for others.

MOMENTUM

SGD has troubles (i.e. it oscillates) in areas where the surface curves much more steeply in one dimension than in another (which are common around local optima).

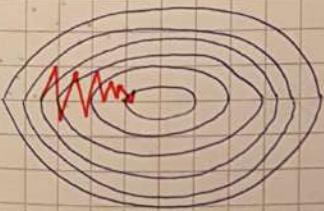
→ Therefore, we consider the MOMENTUM: in that way, the update is the linear combination of previous gradient and new one.

The momentum parameter γ is usually set to 0.9 or a similar value

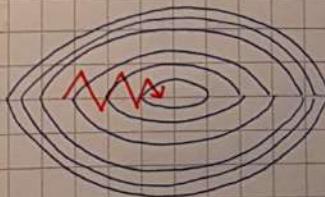
$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

It helps accelerate SGD in relevant direction and dampens oscillations.

Remark - The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.



SGD without momentum



SGD with momentum

ADVANCED SGD SCHEMES & algorithms

- ADAGRAD adapts the learning rate for each parameter independently
 - it performs smaller updates (i.e. low learning rates) for parameters associated with ~~infrequent features~~ frequently occurring features.
 - it performs larger updates (i.e. high learning rates) for parameters associated with infrequent features.
- ADA DELTA (improved version of ADAGRAD)

- RMSprop (improved version of ADAGRAD).
- ADAM (Adaptive Moment Estimation)
 - it also computes adaptive learning rates for each parameter.
 - it combines ideas from ADAGRAD and momentum.
 - whereas momentum can be seen as a ball running down a slope, ADAM behaves like a heavy ball with friction, which thus prefers flat minima in the error surface.

SUPPORT VECTOR MACHINES (cap. 15)

Aim: learning linear predictors in high dimensional feature spaces.

The SVM algorithmic paradigm tackles the sample complexity challenge by searching for "large margin" separators. Roughly speaking, a halfspace separates a training set with a large margin if all the examples are not only on the correct side of the separating hyperplane but also far away from it. Restricting the algorithm to output a large margin separator can yield a small sample complexity even if the dimensionality of the feature space is high (even infinite).

CLASSIFICATION MARGIN

Consider a classification problem with two classes:

Training data: $S = ((\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m))$
with $\vec{x}_i \in \mathbb{R}^d$ and $y_i \in \{\pm 1\}$

Hypothesis set: $H = \text{halfspaces}$

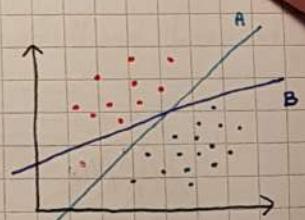
Assumption: the data is linearly separable; there exist a halfspace that perfectly classify the training set.

Find a solution: there are multiple separating hyperplanes that correctly classify the training set, which one is the best?

To answer, we use the concept of "margin".

MARGIN: minimum distance from an example in the training set.

Idea: best separating hyperplane is the one with the largest margin
→ can tolerate more "noise".



LINEARLY SEPARABLE TRAINING SET:

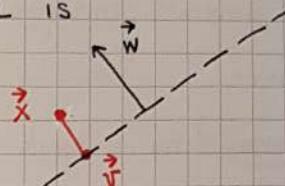
A training set $S = ((\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m))$ is linearly separable if there exists a halfspace (\vec{w}, b) such that $y_i = \text{sigm}(\langle \vec{w}, \vec{x}_i \rangle + b) \forall i=1,\dots,m$

- i.e., it perfectly separates all data in the training set.
- or, equivalently $t_i : y_i(\langle \vec{w}, \vec{x}_i \rangle + b) > 0$

MARGIN:

Given a separating hyperplane defined by $L = \{\vec{v} : \langle \vec{v}, \vec{w} \rangle + b = 0\}$ and given a sample \vec{x} , the distance of \vec{x} to L is

$$d(\vec{x}, L) = \min \{ \|\vec{x} - \vec{v}\| : \vec{v} \in L \}$$



Claim

The distance between a point \vec{x} and the hyperplane defined by (\vec{w}, b) if $\|\vec{w}\| = 1$ is $d(\vec{x}, L) = |\langle \vec{w}, \vec{x} \rangle + b|$

In this case ($\|\vec{w}\|=1$) the margin is:

$$\min_i |\langle \vec{w}, \vec{x}_i \rangle + b|$$

The closest examples are called support vectors.

SUPPORT VECTOR MACHINE (SVM)

HARD-SVM: seek for the separating hyperplane with largest margin.
(NB. works only for linearly separable data)

Computational problem:

$$\arg \max_{(\vec{w}, b)} \min_i |\langle \vec{w}, \vec{x}_i \rangle + b| \text{ st. } \forall i, y_i(\langle \vec{w}, \vec{x}_i \rangle + b) > 0 \quad (\text{linearly separable})$$

equivalent formulation: $\arg \max_{(\vec{w}, b)} \min_{i \in [m]} y_i(\langle \vec{w}, \vec{x}_i \rangle + b) \quad (1)$

An equivalent formulation of the HARD-SVM rule is expressing it as a quadratic formulation problem:

ALGORITHM: HARD-SVM

INPUT: $S = ((\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m))$

SOLVE:

$$(\vec{w}_0, b_0) = \underset{(\vec{w}, b)}{\operatorname{argmin}} \|\vec{w}\|^2 \text{ s.t. } \forall i, y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 \quad (2)$$

$$\text{OUTPUT: } \hat{\vec{w}} = \frac{\vec{w}_0}{\|\vec{w}_0\|}, \quad \hat{b} = \frac{b_0}{\|\vec{w}_0\|}$$

The objective of that algorithm is a convex quadratic function, constraints are linear inequalities: can be solved with quadratic programming solvers.

Lemma The output of Hard-SVM is a solution of equation (1).

Proof. Let (\vec{w}^*, b^*) be a solution of (1) and define the margin achieved by (\vec{w}^*, b^*) to be $\gamma^* = \min_{i \in [m]} y_i (\langle \vec{w}^*, \vec{x}_i \rangle + b^*)$.

Therefore, for all i we have:

$$y_i (\langle \vec{w}^*, \vec{x}_i \rangle + b^*) \geq \gamma^*$$

or equivalently

$$y_i (\langle \frac{\vec{w}^*}{\gamma^*}, \vec{x}_i \rangle + \frac{b^*}{\gamma^*}) \geq 1$$

Hence, the pair $(\frac{\vec{w}^*}{\gamma^*}, \frac{b^*}{\gamma^*})$ satisfies the conditions of the quadratic optimization problem given in equation (2). Therefore, $\|\vec{w}_0\| \leq \|\frac{\vec{w}^*}{\gamma^*}\| = \frac{1}{\gamma^*}$. It follows that for all i ,

$$y_i (\langle \hat{\vec{w}}, \vec{x}_i \rangle + \hat{b}) = \frac{1}{\|\vec{w}_0\|} y_i (\langle \vec{w}_0, \vec{x}_i \rangle + b_0) \geq \frac{1}{\|\vec{w}_0\|} \geq \gamma^*$$

Since $\|\hat{\vec{w}}\| = 1$ we obtain that $(\hat{\vec{w}}, \hat{b})$ is an optimal solution of (1).

HOMOGENEOUS REPRESENTATION

Formulation with homogeneous halfspaces: \rightarrow pass through the origin and are thus defined by $\operatorname{sign}(\langle \vec{w}, \vec{x} \rangle)$, where the bias term b is set to zero

- Assume first component of $x \in X$ is 1, then

$$\vec{w}_0 = \underset{\vec{w}}{\operatorname{argmin}} \|\vec{w}\|^2 \text{ s.t. } \forall i, y_i \langle \vec{w}, \vec{x}_i \rangle \geq 1 \quad (3) \text{ (HARD-SVM)}$$

- Notice that is similar but not exactly the same as the non-homogeneous one (the bias now also goes inside the regularization)
 \rightarrow in practice, no big difference.

case in which we reduce non homogeneous halfspaces by increasing the dimension to d+1 (homogeneous one)

SUPPORT VECTORS

The SUPPORT VECTORS are the vectors at minimum distance from \vec{w}_0 . They are the only training vectors that matter for defining \vec{w}_0 !

(2)

Theorem

Let \vec{w}_0 be as defined in (3) and let $I = \{i : |\langle \vec{w}_0, \vec{x}_i \rangle| = 1\}$.

Then, there exist coefficients a_1, \dots, a_m such that

$$\vec{w}_0 = \sum_{i \in I} a_i \vec{x}_i$$

The examples $\{\vec{x}_i : i \in I\}$ are called "support vectors".

Remark. Solving HARD-SVM is equivalent to find a_i for the support vectors ($a_i \neq 0$ only for support vectors).

DUALITY

Consider the function

$$g(\vec{w}) = \max_{\vec{z} \in \mathbb{R}^m : \vec{z} \geq 0} \sum_{i=1}^m a_i (1 - y_i \langle \vec{w}, \vec{x}_i \rangle) = \begin{cases} 0 & \text{if } \forall i, y_i \langle \vec{w}, \vec{x}_i \rangle \geq 1 \\ \infty & \text{otherwise} \end{cases}$$

We can therefore rewrite equation (3) as

$$\min_{\vec{w}} (\|\vec{w}\|^2 + g(\vec{w}))$$

$$= \frac{1}{2} \cdot \min_{\vec{w}} \max_{\vec{z} \in \mathbb{R}^m : \vec{z} \geq 0} \left(\frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^m a_i (1 - y_i \langle \vec{w}, \vec{x}_i \rangle) \right)$$

It's called the "dual" problem.

Key property: only requires the inner product between instances but not the direct access to instances. It will be very useful for the "Kernel trick".

SOFT-SVM

Key issue: hard-SVM needs the data to be linearly separable \rightarrow almost never true in practical problems.

We need an approach that can work also with non linearly separable data \rightarrow SOFT-SVM.

Relax the constraints of HARD-SVM but take into account the violations of the separation into the objective function.

A natural relaxation is to allow the constraint to be violated for some of the examples in the training set.

- Relax the constraint:

- introduce slack variables : $\vec{z} = (z_1, \dots, z_m)$, $z_i \geq 0$

- for each $i = 1, \dots, m$: $y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - z_i$

- z_i : how much the constraint is violated.

- Soft-SVM jointly minimizes:

- the norm of \vec{w} (\rightarrow maximize margin)

- the average of z_i (\rightarrow minimize constraint violation from misclassified points)

The trade-off between the two objectives is controlled by a parameter $\lambda > 0$. This leads to the SOFT-SVM optimization problem:

Input: $(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)$, parameter $\lambda > 0$

Solve:

$$\min_{\vec{w}, b, \vec{z}} \left(\lambda \|\vec{w}\|^2 + \frac{1}{m} \sum_{i=1}^m z_i \right)$$

$$\text{s.t. } \forall i, y_i (\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - z_i \text{ and } z_i \geq 0$$

Output: \vec{w}, b

REFORMULATE WITH HINGE LOSS

HINGE LOSS:

$$l^{\text{hinge}}((\vec{w}, b), (\vec{x}, y)) = \max \{0, 1 - y (\langle \vec{w}, \vec{x} \rangle + b)\}$$

The problem can be reformulated with the Hinge loss:

$$\min_{\vec{w}, b} \left(\lambda \|\vec{w}\|^2 + \underbrace{\frac{1}{m} \sum_{i=1}^m l^{\text{hinge}}((\vec{w}, b), (\vec{x}_i, y_i))}_{L_s^{\text{hinge}}(\vec{w}, b)} \right)$$

It is often more convenient to consider SOFT-SVM for learning a homogeneous halfspace, where the bias term b is set to be zero, which yields the following optimization problem:

$$\min_{\vec{w}} \left(\lambda \|\vec{w}\|^2 + L_s^{\text{hinge}}(\vec{w}) \right)$$

where

$L_s^{\text{hinge}}(\vec{w})$

Approaches

\rightarrow use S

\rightarrow use S

SGD FOR

HINGE LO

$f^{\text{hinge}}(\vec{w})$

Subgradi

\vec{v}

Update \vec{w}

\vec{w}

\uparrow

stam

We wan

$\min_{\vec{w}}$

\vec{w}

Rum

$$L_s^{\text{hinge}}(\vec{w}) = \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y_i \langle \vec{w}, \vec{x}_i \rangle\}$$

Approaches to solve the problem:

- use standard solvers for optimization problems
- use stochastic gradient descent (SGD)

SGD FOR SOFT SVM

HINGE LOSS:

$$f^{\text{hinge}}(\vec{w}) = \max \{0, 1 - y \langle \vec{w}, \vec{x} \rangle\}$$

Subgradient of f^{hinge} at \vec{w} :

$$\vec{v}^{\text{hinge}} = \begin{cases} 0 & \text{if } 1 - y \langle \vec{w}, \vec{x} \rangle \leq 0 \\ -y \vec{x} & \text{if } 1 - y \langle \vec{w}, \vec{x} \rangle > 0 \end{cases}$$

Update rule (for the complete soft-SVM optimization):

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \vec{v}^{(t)} \quad \text{or} \quad \vec{w}^{(t+1)} = -\frac{1}{\lambda t} \sum_{j=1}^t \vec{v}_j$$

↑

standard SGD

↑

variant of SGD for λ -strongly convex functions

We want to solve:

$$\min_{\vec{w}} \left(\frac{\lambda}{2} \|\vec{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y_i \langle \vec{w}, \vec{x}_i \rangle\} \right)$$

Remark. It's standard to add a $1/2$ in the regularization term to simplify some computations.

Algorithm: SGD for solving soft-SVM

(Algorithm 1) see Kernel section

parameter: T

initialize: $\vec{\theta}^{(0)} = 0$

for $t = 1, \dots, T$:

$$\text{let } \vec{w}^{(t)} = \frac{1}{\lambda t} \vec{\theta}^{(t)}$$

choose i uniformly at random from $\{1, \dots, m\}$:

if $y_i \langle \vec{w}^{(t)}, \vec{x}_i \rangle < 1$ then $\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} + y_i \vec{x}_i$

$$\text{else } \vec{\theta}^{(t+1)} = \vec{\theta}^{(t)}$$

$$\text{output: } \vec{w} = \frac{1}{T} \sum_{t=1}^T \vec{w}^{(t)}$$

$$\vec{\theta}^{(t)} = -\sum_{j=1}^{t-1} \vec{v}_j$$

hinge loss: gradient is \emptyset if correctly classified
and $-y_i \vec{x}_i$ if error

KERNEL METHODS (cap. 16)

In the previous chapter we tackle the sample complexity issue using the concept of margin. In this chapter we tackle the computational complexity challenge using the method of Kernels.

SVM is a powerful algorithm, but still limited to linear models, and linear models cannot always be used (directly!).

→ Idea: • apply a nonlinear transformation to each point in training;

- learn a linear predictor in the transformed space;
- make a prediction for a new instance.

We start the chapter by describing the idea of embedding the data into a high dimensional feature space. We then introduce the idea of Kernels. A Kernel is a type of a similarity measure between instances. The special property of Kernel similarities is that they can be viewed as inner products in some Hilbert space to which the instance space is virtually embedded.

EMBEDDINGS INTO FEATURE SPACES

Define a non-linear mapping ψ from the input space to a new (larger) space:

1. Given a domain set X and a learning task, find a mapping to a new feature space F :

$$\psi: X \rightarrow F$$

F is usually \mathbb{R}^m for some m , but can be an arbitrary Hilbert space (even of infinite size).

2. Given a sequence of labeled examples $S = ((\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m))$ map them to $\hat{S} = ((\psi(\vec{x}_1), y_1), \dots, (\psi(\vec{x}_m), y_m))$

3. Train a linear predictor h over \hat{S} .

4. Predict ~~the label~~ the label of \vec{x} as $h(\psi(\vec{x}))$.

The success of this learning paradigm depends on choosing a good ψ for a given learning task: that is, a ψ that will make the image of the data

distribution linearly separable in the feature space, thus making the resulting algorithm a good learner for a given task.

THE KERNEL TRICK

- The learning over the new highly dimensional space makes halfspaces more expressive.

- On the other side, the computational complexity can become huge.

The solution: Kernel-based learning.

- Kernel: inner product in the feature space.

- Kernel Function $K(\vec{x}, \vec{x}') = \langle \psi(\vec{x}), \psi(\vec{x}') \rangle$

- $K(\vec{x}, \vec{x}')$ represents similarity of the samples in a space where the similarities are realized as inner products.

- Key result: machine learning algorithms for halfspaces can be carried out just on the basis of the values of the kernel function without explicitly representing the points in the feature space.

- Sometimes we can compute $K(\vec{x}, \vec{x}')$ without computing $\psi(\vec{x})$ and $\psi(\vec{x}')$.

Consider $n \in \mathbb{R}^d$:

$$\psi(\vec{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, x_1x_3, \dots, x_dx_d)^T$$

example with
2nd degree
polynomial

The dimension of $\psi(\vec{x})$ is $1 + d + d^2$

$$\langle \psi(\vec{x}), \psi(\vec{x}') \rangle = 1 + \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j \quad O(d^2)$$

Note that

$$\sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j = \left(\sum_{i=1}^d x_i x'_i \right) \left(\sum_{j=1}^d x_j x'_j \right) = (\langle \vec{x}, \vec{x}' \rangle)^2$$

Therefore

$$K_\psi(\vec{x}, \vec{x}') = \langle \psi(\vec{x}), \psi(\vec{x}') \rangle = 1 + \langle \vec{x}, \vec{x}' \rangle + (\langle \vec{x}, \vec{x}' \rangle)^2 \quad O(d)$$

Remark. Computing $\psi(\vec{x})$ requires $O(d^2)$ time, computing $K_\psi(\vec{x}, \vec{x}')$ from the last formula requires $O(d)$ time.

When $K_\psi(\vec{x}, \vec{x}')$ is efficiently computable, we don't need to explicitly compute $\psi(\vec{x}) \rightarrow$ Kernel trick

KERNEL TRICK APPLIED TO SVM

Note that all versions of the SVM optimization problem we have derived in the previous chapter are instances of the following general problem:

$$\min_{\vec{w}} f(\langle \vec{w}, \psi(\vec{x}_1) \rangle, \dots, \langle \vec{w}, \psi(\vec{x}_m) \rangle) + R(\|\vec{w}\|) \quad (1)$$

where $F: \mathbb{R}^m \rightarrow \mathbb{R}$ arbitrary function

$R: \mathbb{R}_+ \rightarrow \mathbb{R}$ monotonically non decreasing function

Theorem (Representer theorem)

Assume that ψ is a mapping from \mathcal{X} to a Hilbert space. Then, there exists a vector $\vec{d} \in \mathbb{R}^m$ such that $\vec{w} = \sum_{i=1}^m d_i \psi(\vec{x}_i)$ is an optimal solution of equation (1).

Consequence: we can optimize the problem with respect to the coefficients d_i getting a problem that depends only on $K(\vec{x}, \vec{x}') = \langle \psi(\vec{x}), \psi(\vec{x}') \rangle$ without explicitly computing $\psi(\vec{x})$ or $\psi(\vec{x}')$.

Note that:

$$\langle \vec{w}, \psi(\vec{x}_i) \rangle = \left\langle \sum_j d_j \psi(\vec{x}_j), \psi(\vec{x}_i) \right\rangle = \sum_{j=1}^m d_j \langle \psi(\vec{x}_j), \psi(\vec{x}_i) \rangle$$

similarly,

$$\|\vec{w}\| = \left\langle \sum_j d_j \psi(\vec{x}_j), \sum_i d_i \psi(\vec{x}_i) \right\rangle = \sum_{i,j=1}^m d_i d_j \langle \psi(\vec{x}_i), \psi(\vec{x}_j) \rangle$$

Let $K(\vec{x}, \vec{x}') = \langle \psi(\vec{x}), \psi(\vec{x}') \rangle$ be a function that implements the kernel function with respect to the embedding ψ . Instead of solving equation (1) we can solve the equivalent problem

$$\min_{\vec{d} \in \mathbb{R}^m} f\left(\sum_{j=1}^m d_j K(\vec{x}_j, \vec{x}_1), \dots, \sum_{j=1}^m d_j K(\vec{x}_j, \vec{x}_m)\right) + R\left(\sqrt{\sum_{i,j=1}^m d_i d_j K(\vec{x}_i, \vec{x}_j)}\right)$$

To solve the last equation, we solely need to know the value of the $m \times m$ matrix G such that $G_{i,j} = K(\vec{x}_i, \vec{x}_j)$, which is often called the GRAM MATRIX G .

In particular, specifying the preceding to the soft-SVM problem we can rewrite the problem as

$$\min_{\vec{d} \in \mathbb{R}^m} (\lambda \vec{d}^\top G \vec{d} + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i (\vec{G} \vec{d})_i\})$$

Once we know the coefficients \vec{w} we can calculate the prediction on a new instance by

$$\langle \vec{w}, \vec{\psi}(\vec{x}) \rangle = \sum_{j=1}^m d_j \langle \vec{\psi}(\vec{x}_j), \vec{\psi}(\vec{x}) \rangle = \sum_{j=1}^m d_j K(\vec{x}_j, \vec{x})$$

The advantage of working with kernels rather than directly optimizing \vec{w} in the feature space is that in some situations the dimension of the feature space is extremely large while implementing the kernel function is very simple.

→ POLYNOMIAL KERNELS

$$K(\vec{x}, \vec{x}') = (1 + \langle \vec{x}, \vec{x}' \rangle)^K$$

K is the degree of the polynomial

- it is a Kernel Function, i.e., $K(\vec{x}, \vec{x}') = \langle \vec{\psi}(\vec{x}), \vec{\psi}(\vec{x}') \rangle$
- $\vec{\psi}: \mathbb{R}^m \rightarrow \mathbb{R}^{(m+1)K}$ contains all the monomials up to degree K .
- Halfspace over $\vec{\psi}$ corresponds to a polynomial predictor of ~~multiple~~ order K in the original space.
- Complexity of computation is $O(m)$ while the dimension of feature space is $O(m^K)$.

→ GAUSSIAN KERNELS (radial basis functions, RBF)

$$K(\vec{x}, \vec{x}') = \exp \left\{ -\|\vec{x} - \vec{x}'\|^2 / 2\sigma^2 \right\}$$

- it is a Kernel Function, i.e., $K(\vec{x}, \vec{x}') = \langle \vec{\psi}(\vec{x}), \vec{\psi}(\vec{x}') \rangle$
- The feature space is of infinite dimension, but computing the kernel is simple and fast!
- Product is close to 0 if instances are far and close to 1 if they are close.
- Parameter σ controls what we mean by "close".
- We can learn any polynomial predictor in the original space by using a Gaussian Kernel.
- VC-dimension is infinite (sample complexity depends on the margin in the feature space).

SOFT SVM WITH KERNELS

We want to solve :

$$\min_{\vec{w}} \left(\lambda \|\vec{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - y_i \langle \vec{w}, \vec{\psi}(\vec{x}_i) \rangle\} \right) \quad (1)$$

Let :

- K be the Kernel function : $K(\vec{x}, \vec{x}') = \langle \vec{\psi}(\vec{x}), \vec{\psi}(\vec{x}') \rangle$
- We shall maintain two vectors in \mathbb{R}^m , corresponding to two vectors $\vec{\beta}^{(t)}$ and $\vec{w}^{(t)}$ defined in the SGD procedure (algorithm 1). That is $\vec{\beta}^{(t)}$ will be a vector such that :

$$\vec{\beta}^{(t)} = \sum_{j=1}^m \beta_j^{(t)} \vec{\psi}(\vec{x}_j),$$

and $\vec{d}^{(t)}$ be such that

$$\vec{w}^{(t)} = \sum_{j=1}^m d_j^{(t)} \vec{\psi}(\vec{x}_j)$$

The vectors $\vec{\beta}$ and \vec{d} are updated according to the following procedure.

Algorithm: SGD for solving SOFT-SVM with Kernels

Goal: solve equation (1)

parameter: T

initialize : $\vec{\beta}^{(0)} = 0$

for $t = 1, \dots, T$

$$\text{let } \vec{d}^{(t)} = \frac{1}{\lambda t} \vec{\beta}^{(t)}$$

choose i uniformly at random from $[m]$

For all $j \neq i$ set $\beta_j^{(t+1)} = \beta_j^{(t)}$

If $(y_i \sum_{j=1}^m d_j^{(t)} K(\vec{x}_j, \vec{x}_i) < 1)$

$$\text{set } \beta_i^{(t+1)} = \beta_i^{(t)} + y_i$$

else

$$\text{set } \beta_i^{(t+1)} = \beta_i^{(t)}$$

output : $\vec{w} = \sum_{i=1}^m d_i \vec{\psi}(\vec{x}_i)$ where $\vec{d} = \frac{1}{T} \sum_{i=1}^T \vec{d}^{(t)}$