# PAC learning model
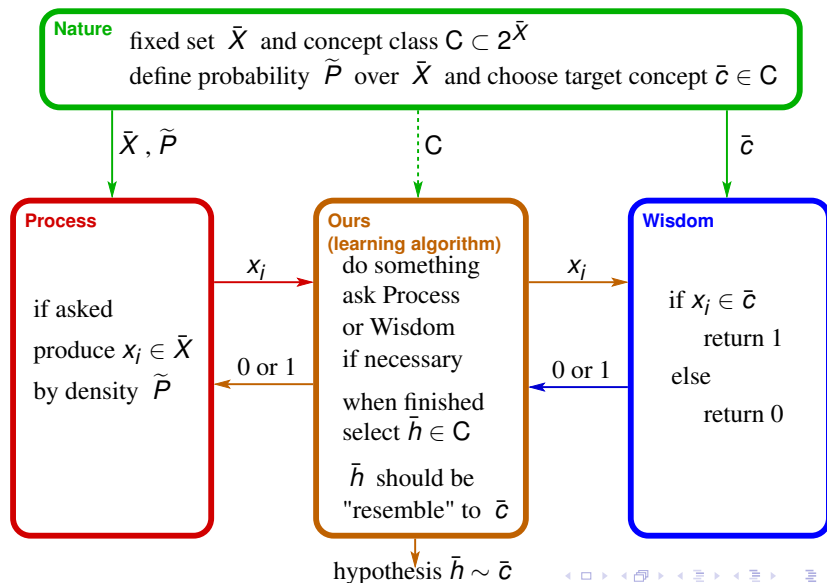
## František Hakl

ICS AS CR
hakl@cs.cas.cz

Mar 2013

" $\bar{h}$ should be resemble to $\bar{c}$ "



### Definition

1. $e_{\widetilde{P}}\left(\bar{h}, \bar{c}\right) \stackrel{\text{def}}{=} \widetilde{P}\left(\bar{c} \triangle \bar{h}\right)\left(=\left(\bar{c} \dot{-} \bar{h}\right) \cup \left(\bar{h} \dot{-} \bar{c}\right)\right)$

2. $\bar{h}$ is consistent if and only if $\{x_i, \ldots, x_m\} \cap \left(\bar{c} \triangle \bar{h}\right) = \emptyset$

### Definition (sample space)

Let $\breve{x} \stackrel{\text{def}}{=} \{x_1, \ldots, x_m\}$, $x_i \in \bar{X}$, $i \in \{1, \ldots, m\}$, $\vec{z} \in \{-1, +1\}^m$ and let $\bar{c} \subset \bar{X}$. Then the ordered tuple

$$\left( \breve{x}, \vec{z} \right)$$

is a $m$-SAMPLE OF CONCEPT $\bar{c}$ if and only if

$$(\forall i \in \{1, \ldots, m\}) \left( (x_i \in \bar{c}) \Leftrightarrow (\vec{z}_i = 1) \right).$$

For concept class C define SAMPLE SPACE OF CONCEPT CLASS as

$$\bar{S}_C \stackrel{\text{def}}{=} \bigcup_{m \geq 1} \left\{ \bigcup_{\bar{c} \in C} \left\{ \left( \breve{x}, \vec{z} \right) \middle| \left( \breve{x}, \vec{z} \right) \text{ is a } m\text{-sample of concept } \bar{c} \right\} \right\}.$$

### Definition $((\epsilon, \delta)$-learning algorithm)

1. $(\epsilon, \delta)$-LEARNING ALGORITHM is each mapping $\widetilde{A^*} : \bar{S}_C \to C$ such that for all $\bar{c} \in C$, $\epsilon, \delta \in (0, 1)$ and $\widetilde{P}$ on $\bar{X}$, the probability of the set

$$\left\{ \widecheck{x} \,\middle|\, \left(\widecheck{x}, \vec{z}\right) \text{ is } m\text{-sample of } \bar{c} \text{ and } e_{\widetilde{P}}\left(\bar{c}, \widetilde{A^*}\left(\left(\widecheck{x}, \vec{z}\right)\right)\right) \geq \epsilon \right\}$$

is smaller than the number $\delta$.

2. If such a learning algorithm exists we say that C IS UNIFORMLY LEARNABLE .

### Theorem

*Let us assume that* $C$ *is a concept class over finite set* $\bar{X}$ *and* $H = C$. *Then, for each learning algorithm* $A^*$ *requiring*

$$\frac{1}{\epsilon} \ln \left( \frac{|C|}{\delta} \right)$$

*queries and producing for the given concept* $\bar{c} \in C$ *a consistent hypothesis it holds that*

$$Prob_{\widetilde{P}} \left( e_{\widetilde{P}} \left( \bar{c}, \widetilde{A^*} \left( \left( \overset{\smile}{x}, \bar{z} \right) \right) \right) \geq \epsilon \right) < \delta.$$

### Definition

For any arbitrary $R \subset 2^{\bar{X}}$ and for any arbitrary probability density $\widetilde{P}$ defined on $\bar{X}$ and for an arbitrary $\epsilon > 0$ let us define $R_{\widetilde{P}, \epsilon} \overset{\text{def}}{=} \{\bar{r} \in R \,|\, Prob_{\widetilde{P}}(\bar{r}) > \epsilon\}$. Then, we will call $\bar{T}_{\widetilde{P}, \epsilon} \subset \bar{X}$ $\epsilon$-TRANSVERSAL $R$ just if

$$\left(\forall \bar{r} \in R_{\widetilde{P}, \epsilon}\right) \left(\bar{r} \cap \bar{T}_{\widetilde{P}, \epsilon} \neq \emptyset\right)$$

### Example

$\bar{X} = \langle 0, 1 \rangle^n$, $\widetilde{P}$ uniform on $\bar{X}$, $R = \{\bar{b} \subset \bar{X} \,|\, \bar{b} \text{ is a ball}\}$, $\epsilon = \frac{1}{z}$. Then $\bar{T} = \{k\epsilon \,|\, k = 0, \cdots, \frac{1}{\epsilon}\}^n$ is a $\left(\frac{\pi^{\frac{n}{2}}(\sqrt{n}\epsilon)^n}{\Gamma(\frac{n}{2}+1)2^n}\right)$-transversal of $R$.

... if hypotheses $\bar{h}$ produced by an algoritm is consistent and has $e_{\widetilde{P}}(\bar{h}, \bar{c}) > \epsilon$, then $\{x_i, \ldots, x_m\}$ can't be $\epsilon$-transversal of the system $R \stackrel{\text{def}}{=} \{\bar{h} \triangle \bar{c} \mid \bar{h} \in H\}$ ...

### Definition

For each $m \geq 1$, $\epsilon > 0$ let

$$\bar{Q}_{m,\epsilon} \stackrel{\text{def}}{=} \left\{ \breve{x} \in \bar{X}^m \, \middle| \, \breve{x} \text{ do not form } \epsilon\text{-transversal of } R \right\}$$

and (assume that $\breve{x}, \breve{y} \in \bar{X}^m$)

$$\bar{J}_\epsilon^{2m} \stackrel{\text{def}}{=} \left\{ \breve{xy} \in \bar{X}^{2m} \, \middle| \, \left( \exists \bar{r} \in R_{\widetilde{P},\epsilon} \right) \left( \breve{x} \cap \bar{r} = \emptyset \text{ and } \left| \breve{y} \cap \bar{r} \right| \geq \frac{\epsilon m}{2} \right) \right\}.$$

... the probability of the set $\bar{Q}_{m,\epsilon}$ is a probability of producing consistent hypothesis with error $e_{\widetilde{P}}(\bar{h}, \bar{c}) > \epsilon$ ...

### Definition

1. The class  H  is well-behaved if the sets  $\bar{Q}_{m,\epsilon}$  and  $\bar{J}_\epsilon^{2m}$  are measurable for any probability  $\widetilde{P}$ , any  $m \geq 1$ ,  $\epsilon > 0$ , and any system of sets  $R \stackrel{\text{def}}{=} \{\bar{h} \triangle \bar{c} \,|\, \bar{h} \in H\}$ , where  $\bar{c}$  is an arbitrary Borelian set.

2. The class  $H \subset 2^{\bar{X}}$  is universally separable, if there exists a countable subset  T  of the class  H  such that for all  $\bar{h} \in H$  there exists a sequence  $\{\bar{h}_i\}_1^\infty$  of sets from  T  such that

$$\left(\forall x \in \bar{X}\right)\left(\exists n \geq 1\right)\left(\left(\forall i \geq n\right)\left(x \in \bar{h}_i \text{ if and only if } x \in \bar{h}\right)\right).$$

### Theorem

*If*  H  *is universally separable, then*  H  *is well-behaved.*

### Lemma

*Let* $R \neq \emptyset$ *be a concept class and* $\widetilde{P}$ *be probability on* $\bar{X}$ *for which* $\bar{Q}_{m,\epsilon}$ *and* $\bar{J}_\epsilon^{2m}$ *are measurable for all* $m \geq 1$, $\epsilon > 0$. *Then for each* $\epsilon > 0$ *and* $m \geq \frac{2}{\epsilon}$

$$Prob_{\widetilde{P}^m}\left(\bar{Q}_{m,\epsilon}\right) < 2Prob_{\widetilde{P}^{2m}}\left(\bar{J}_\epsilon^{2m}\right) \leq 2\Pi_R\left(2m\right)2^{-\frac{\epsilon m}{2}}.$$

### Lemma

*Let* $m \geq 1$, $\bar{c} \subset \bar{X}$, $H \subset 2^{\bar{X}}$ *and* $R \stackrel{def}{=} \left\{\bar{h} \bigtriangleup \bar{c} \,\middle|\, \bar{h} \in H\right\}$. *Then*

$$\Pi_R\left(m\right) = \Pi_H\left(m\right).$$

### Lemma

*Let* $d = \text{VC}_{dim}\left(H\right) \in N$, $\epsilon, \delta \in (0, 1)$, $\Gamma \stackrel{def}{=} 2\Pi_H\left(2m\right)2^{-\frac{\epsilon m}{2}}$. *Then*
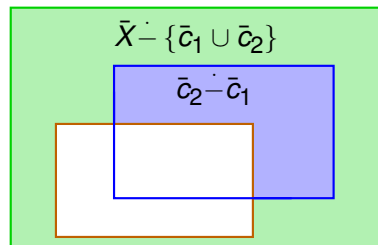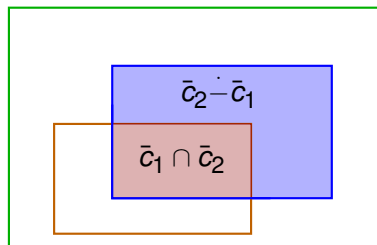
$$m \geq \max\left\{\frac{4}{\epsilon}\log_2\left(\frac{2}{\delta}\right), \frac{8d}{\epsilon}\log_2\left(\frac{12.611}{\epsilon}\right)\right\} \Rightarrow \Gamma \leq \delta.$$

### Definition

Concept class C is called nontrivial iff

$$(\exists \bar{c}_1, \bar{c}_2 \in \mathsf{C})\,(\bar{c}_1 \neq \bar{c}_2 \text{ and } (\bar{c}_1 \cap \bar{c}_2 \neq \emptyset \text{ or } \bar{c}_1 \cup \bar{c}_2 \neq \bar{X}))\ .$$

Concept class is called trivial in other cases.



Two cases of minimal content of nontrivial concept class.
(colored sets are nonempty)

## Theorem (main result of PAC theory)

*Let* C *be a nontrivial, well-behaved class. Then:*

1. *If* $VC_{dim}(C) = d < +\infty$. *Then*
   1. *for any* $0 < \epsilon < \frac{1}{2}$ *there is no* $(\epsilon, \delta)$*-learning algorithm with number of queries less than*

   $$\max\left(\frac{1-\epsilon}{\epsilon}\ln\left(\frac{1}{\delta}\right), d\left(1 - 2\left(\epsilon\left(1 - \delta\right) + \delta\right)\right)\right) . \quad (1)$$

   2. *for arbitrary* $0 < \epsilon < 1$*, any learning algorithm using at least*

   $$\max\left(\frac{4}{\epsilon}\log_2\left(\frac{2}{\delta}\right), \frac{8d}{\epsilon}\log_2\left(\frac{12.611}{\epsilon}\right)\right) \quad (2)$$

   *queries and returning a consistent hypothesis is an* $(\epsilon, \delta)$*-learning algorithm.*

2. C *is uniformly learnable* **if and only if** $VC_{dim}(C) < +\infty$.

Sketch of the proof:

1. 1.
   - $\frac{1-\epsilon}{\epsilon} \ln\left(\frac{1}{\delta}\right)$: (c&c) Any nontrivial concept class can be reduced to one of the cases discussed above. For uniform probability we get a contradiction.
   - $d\left(1 - 2\left(\epsilon\left(1 - \delta\right) + \delta\right)\right)$: (c&c) Reduce $\bar{X}$ to $d$-element subset with uniform probability. Then use the "matrix" $\boldsymbol{Z}_{\bar{c},\bar{h}} \stackrel{\text{def}}{=} e_{\widetilde{P}}\left(\bar{c}, \bar{h}\right)$ to show, that $m > d\left(1 - 2\left(\epsilon\left(1 - \delta\right) + \delta\right)\right)$ imply that $\left(\exists \bar{h}^*\right)$ contradicts $\left(\epsilon, \delta\right)$-property . . . "broadly speaking".

   2. See previous slides.

2. 
   - $\Leftarrow$ (construction) Use Zermelo's well-ordering theorem to well-order $\bar{H}$. Let algorithm get $m$-sample of $\bar{c}$ and return the first hypothesis consistent with $\bar{c}$. The statemet follows from 1)-2).
   - $\Rightarrow$ (by contradiction) For any $d \in N$ we carry out steps 1)-1)-(second term). Choose $\left(\epsilon, \delta\right)$ such that $\left(1 - 2\left(\epsilon\left(1 - \delta\right) + \delta\right)\right) > 0$. Hence $m$ can't be upper-bounded.

František Hakl                    ICS AS CR

### Definition (discrete delta rule)

Let $\left(\vec{x}_1, y_1\right), \ldots, \left(\vec{x}_t, y_t\right)$ be a given sequence of tuples in $\Re^n \times \{-1, +1\}$, $t \geq 1$. Further, let vector's sequence $\{\vec{w}_i\}_1^\infty$ satisfy the following recursive formulas

1. put $\vec{w}_1 \stackrel{\text{def}}{=} \vec{0}$, $k = 1$

2. let $k = k + 1$ and $\bar{J} \stackrel{\text{def}}{=} \{j \in \{1, \ldots, t\} \mid \widetilde{sgn}\left(\langle \vec{w}_k \mid \vec{x}_j \rangle\right) \neq y_j\}$

   1. if $\bar{J} = \emptyset$ put $\vec{w}_{k+1} = \vec{w}_k$ and STOP,
   2. else let $j_k \in \bar{J}$ be arbitrary. Then put

   $$\vec{w}_{k+1} \stackrel{\text{def}}{=} \vec{w}_k + y_{j_k} \vec{x}_{j_k}$$

   and REPEAT step 2).

Then we say that $\{\vec{w}_i\}_1^\infty$ is DELTA SEQUENCE of $\left(\vec{x}_1, y_1\right), \ldots, \left(\vec{x}_t, y_t\right)$.

### Theorem (delta rule convergence)

*Let $\{\vec{w}_i\}_1^\infty$ is a delta sequence and let there exists a vector $\widehat{\widehat{\boldsymbol{w}}}$
such that for all indexes $i \in \{1, \ldots, t\}$ holds
$\widetilde{sgn}\left(\left\langle \widehat{\widehat{\boldsymbol{w}}} \mid \vec{\boldsymbol{x}}_i \right\rangle\right) = y_i$. Further let*

$$\alpha \stackrel{def}{=} \max_{i \in \{1, \ldots, t\}} \left\{ \|\vec{\boldsymbol{x}}_i\|^2 \right\} \quad and \quad \beta \stackrel{def}{=} \min_{i \in \{1, \ldots, t\}} \left\{ \left| \left\langle \widehat{\widehat{\boldsymbol{w}}} \mid \vec{\boldsymbol{x}}_i \right\rangle \right| \right\} > 0 \; .$$

*Then there exists an natural number $z > 0$ satisfying
$\vec{w}_{z+1} = \vec{w}_z$ and $z$ can be estimated as*

$$z \leq \frac{\alpha \left\| \widehat{\widehat{\boldsymbol{w}}} \right\|^2}{\beta^2} + 1 \; .$$

### Theorem (delta rule complexity)

*There exists a linearly separable dichotomy of the $\{-1, +1\}^n$ such that any integer linear separator $(\vec{w}, t)$ of this dichotomy satisfies estimation*

$$2^{\frac{n-2}{2}} \leq \sum_{k=1}^{n} \left| \vec{w}_k \right| + |t|.$$

### Definition (Mangasarian LP)

Let $\bar{A} \stackrel{\text{def}}{=} \{\vec{a}_1, \ldots, \vec{a}_i\}$ and $\bar{B} \stackrel{\text{def}}{=} \{\vec{b}_1, \ldots, \vec{b}_j\}$ be a finite subsets of the $\Re^n$. Then MANGASARIAN LINEAR PROBLEM is defined as the problem find vectors $\vec{y} \in \Re^i$, $\vec{z} \in \Re^j$, $\vec{w} \in \Re^n$ and $t \in \Re$ that minimizes

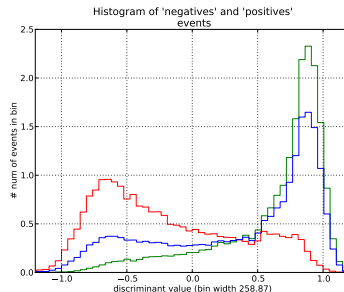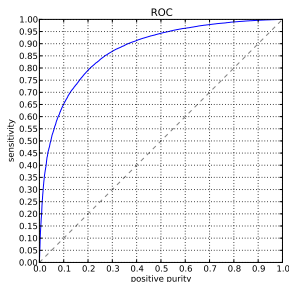$$\sum_{\alpha=1}^{i} \vec{y}_\alpha + \sum_{\beta=1}^{j} \vec{z}_\beta$$

subject to

$$
\begin{array}{rcll}
\vec{y}_\alpha + \left\langle \vec{w} \,\middle|\, \vec{a}_\alpha \right\rangle - t & \geq & 1 & \text{for} \quad \alpha \in \{1, \ldots, i\} \\
\vec{z}_\beta - \left\langle \vec{w} \,\middle|\, \vec{b}_\beta \right\rangle + t & \geq & 1 & \text{for} \quad \beta \in \{1, \ldots, j\} \\
\vec{y}_\alpha & \geq & 0 & \text{for} \quad \alpha \in \{1, \ldots, i\} \\
\vec{z}_\beta & \geq & 0 & \text{for} \quad \beta \in \{1, \ldots, j\} .
\end{array}
$$

### Theorem

Let $\bar{A} \stackrel{def}{=} \{\vec{a}_1, \ldots, \vec{a}_i\}$ and $\bar{B} \stackrel{def}{=} \{\vec{b}_1, \ldots, \vec{b}_j\}$ be a finite subsets of the $\Re^n$ . Then

1. There exists a linear separator of the sets $\bar{A}$ and $\bar{B}$ if and only if the optimal value of the corresponding Mangasarian LP is zero.

2. If the optimal value of the corresponding Mangasarian LP is zero and $(\vec{y}^*, \vec{z}^*, \vec{w}^*, t^*)$ is optimal solution, than $(\vec{w}^*, t^*)$ is linear separator of the sets $\bar{A}$ and $\bar{B}$ .

# . . . a real case . . .



The necessary condition on consistency of hypotheses produced is unsatisfied, PAC model isn't applicable. We have to use

Probably Approximately Optimal (PAO) model

$$\mathcal{P} \; \mathcal{A} \; \mathcal{C}$$