# Estimates of the learning set size for k-NN and IINC methods in HEP
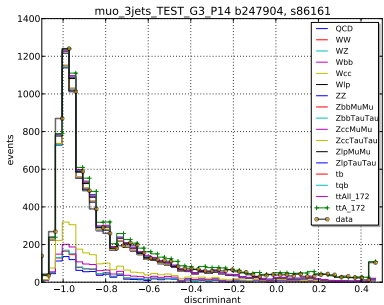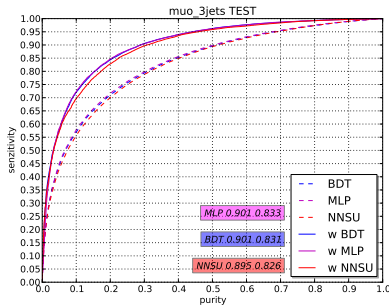
## František Hakl

SPMS 2018

hakl@cs.cas.cz

Institute of computer science, Prague

Jun 2018

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

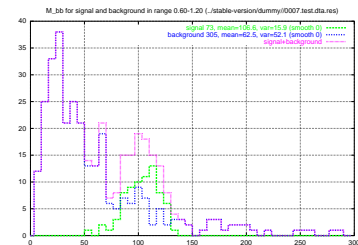k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion
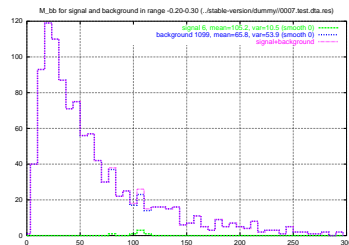
Data driven verification use data to test

- over-learning of training method, e.g. resemble behavior on learn, train and test data
- results robustness on different portions of data (cross-validation)
- resemblance of discriminant distribution over different data sets or over simulated and measured data (so called control plots)

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

- validity of apriori known statistical characteristics of data
  (enhanced by cross-validation)



These approaches do not provide an estimation of
convenient size of data sets and information about expected
accuracy of separation.

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

Probably approximately correct learning

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

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

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

3. $\bar{S}_C$ denote the set of all samples $\left(\breve{x}, \vec{z}\right)$ of fixed $\bar{c} \in C$, where
   $\vec{z} \in \{-1, +1\}^m$, $\breve{x} \in \bar{X}^m$, $m \in Z$.

4. $(\epsilon, \delta)$-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

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

   is smaller than the number $\delta$.

5. VC-dimension: Let $\bar{X}$ be arbitrary set, $C \subset 2^{\bar{X}}$ and

   $$\Pi_C(m) \overset{\text{def}}{=} \max_{\bar{A} \subset \bar{X}, |\bar{A}| = m} \left|\{\bar{b} \,|\, (\exists \bar{c} \in C)\,(\bar{b} = \bar{A} \cap \bar{c})\}\right|$$

   Then

   $$VC_{dim}(C) \overset{\text{def}}{=} \sup\left\{m \,\middle|\, \Pi_C(m) = 2^m\right\}.$$

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
**Main theorem**
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

## Theorem (main result of PAC theory)

*Let* C *satisfy* $(\exists \bar{c}_1, \bar{c}_2 \in C) (\bar{c}_1 \neq \bar{c}_2$ *and* $(\bar{c}_1 \cap \bar{c}_2 \neq \emptyset$ *or* $\bar{c}_1 \cup \bar{c}_2 \neq \bar{X}))$ *and* C *be well-behaved . Then:*

1. *If* $\mathrm{VC}_{dim}(C) < +\infty$. *Then*

   1. *for any* $0 < \epsilon < \frac{1}{2}$ *there is no* $(\epsilon, \delta)$*-learning algorithm with number of queries less than*

   $$\max\left(\frac{1-\epsilon}{\epsilon}\ln\left(\frac{1}{\delta}\right), \mathrm{VC}_{dim}(C) \cdot (1 - 2(\epsilon(1-\delta) + \delta))\right) . \quad (1)$$

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

   $$\max\left(\frac{4}{\epsilon}\log_2\left(\frac{2}{\delta}\right), \frac{8\mathrm{VC}_{dim}(C)}{\epsilon}\log_2\left(\frac{13}{\epsilon}\right)\right) \quad (2)$$

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

2. $(\epsilon, \delta)$*-learning algorithm for* C *exists* $\Leftrightarrow \mathrm{VC}_{dim}(C) < +\infty$.

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

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" $\boldsymbol{Z}_{\bar{c},\bar{h}} \overset{\text{def}}{=} e_{\widetilde{P}}\left(\bar{c},\bar{h}\right)$ to show, that $m > d\left(1 - 2\left(\epsilon\left(1-\delta\right) + \delta\right)\right)$ imply that $(\exists h^*)$ contradicts $(\epsilon, \delta)$-property . . ."broadly speaking".

   2. In more steps we show that from (2) follows that

   $$Prob_{\widetilde{P}}\left(\{x_i, \ldots, x_m\} \mid \left(\forall \bar{T} \in \{\bar{h} \triangle \bar{c} \mid \bar{h} \in H\} \mid Prob_{\widetilde{P}}\left(\bar{T}\right) > \epsilon\right)\right.$$

   $$\left.\left(\{x_i, \ldots, x_m\} \cap \bar{T} = \emptyset\right)\right) \leq \delta .$$

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 statement follows from 1)-2).
   - $\Rightarrow$ (by contradiction) For any $d \in N$ we carry out steps 1)-1)-(second term). Choose $(\epsilon, \delta)$ such that $(1 - 2\left(\epsilon\left(1-\delta\right) + \delta\right)) > 0$. Hence $m$ can't be upper-bounded.

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

Nearest neighbor (NN) is consistent and has a known $\text{VC}_{dim}(NN)$ u.b.

## Lemma (Union, Intersection)

Let $\mathsf{U}_{k,\mathsf{C}} \overset{def}{=} \left\{ \bigcup_{i=1}^{k} \bar{c}_i \,\middle|\, \left(\forall i \in \hat{k}\right)(\bar{c}_i \in \mathsf{C}) \right\}$, $\mathsf{I}_{k,\mathsf{C}} \overset{def}{=} \left\{ \bigcap_{i=1}^{k} \bar{c}_i \,\middle|\, \left(\forall i \in \hat{k}\right)(\bar{c}_i \in \mathsf{C}) \right\}$
and $\text{VC}_{dim}(\mathsf{C}) = d \geq 1$ be finite. Then

$$\text{VC}_{dim}(\mathsf{U}_{k,\mathsf{C}}) \leq 2dk\log_2(3k) \quad and \quad \text{VC}_{dim}(\mathsf{I}_{k,\mathsf{C}}) \leq 2dk\log_2(3k).$$

$\bar{X} = \Re^n$, $k$=number of Balls (or Rect.), $\text{VC}_{dim}(Ball_n) = n + 1$, $\text{VC}_{dim}(Rect_n) = 2n$

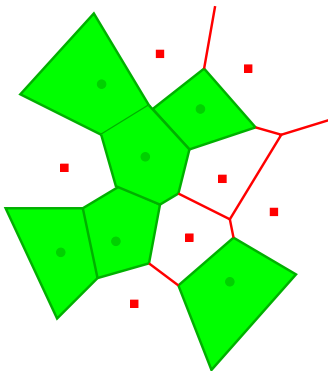**Euclidean**                                        **Manhattan**



$\text{VC}_{dim}(NN_{Ball_n}) \leq 2(n+1)k\log_2(3k)$, consistent, $\text{VC}_{dim}(NN_{Rect_n}) \leq 4nk\log_2(3k)$

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

IINC algorithm is consistent and has a known $\mathrm{VC}_{dim}\left(IINC\right)$ upper bound

## Lemma

*Let $\bar{X}$ be an arbitrary set, $\mathrm{C} \subset 2^X$. Then*

1. *if any two sets in $\mathrm{C}$ are disjoint then $\mathrm{VC}_{dim}\left(\mathrm{C}\right) = 1$,*
2. $\mathrm{VC}_{dim}\left(\mathrm{C}\right) = 1 \Rightarrow \mathrm{VC}_{dim}\left(\mathrm{U}_{k,\mathrm{C}}\right) \leq k.$



IINC outline (basic)

- for a new (unspecified) point compute distances to all *k* known points
- sort points by inverted distances
- put new point to the set of the first point in sorted sequence
- ... so the "green" set is an union of pairwise disjoint sets

$\mathrm{VC}_{dim}\left(IINC\right) \leq k$ and consistent hypothesis

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

## Corollary

Let $|\{x_i, \ldots, x_m\} \cap \bar{c}| = \rho m$, e.g. $\rho$ is the ratio of positive examples. It follows (recall second lower bound $m > \frac{8 \text{VC}_{dim}(C)}{\epsilon} \log_2 \left( \frac{13}{\epsilon} \right)$, $n$ is dimension of examples):

NN Euclidean: $\quad 1 > \rho \cdot 16(n+1) \quad \times \quad \log_2 (3\rho m) \times \left[ \frac{1}{\epsilon} \log_2 \left( \frac{13}{\epsilon} \right) \right]$

NN Manhattan: $\quad 1 > \rho \cdot 32n \quad \times \quad \log_2 (3\rho m) \times \left[ \frac{1}{\epsilon} \log_2 \left( \frac{13}{\epsilon} \right) \right]$

IINC: $\quad 1 > \rho \cdot 8 \quad \times \quad \left[ \frac{1}{\epsilon} \log_2 \left( \frac{13}{\epsilon} \right) \right]$

(note that the first constrain implies $\log_2 (3\rho m) > \log_2 \left( \frac{12\rho}{\epsilon} \log_2 \left( \frac{2}{\delta} \right) \right)$)

## Discussion

- unusable for "large" values of $\rho$ (e.g. $\rho \simeq \epsilon$)
- dimension of examples can be considered constant; corresponds to the number of relevant and reasonable features
- for NN $\rho$ should be proportionate to $\quad \dfrac{\text{desired accuracy of separation } (\epsilon)}{\text{logarithm of positive examples}}$
- applicable in the case of very rare positive examples

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

Example of HEP data set size
(source Measurement of Electroweak Top Quark Production at DØ, Yun-Tse Tsai,
Rochester, 2013)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Pre-tagged event yields** | | | | | | | | |
| | Run IIa, 1 fb$^{-1}$ | | | | Run IIb, 8.7 fb$^{-1}$ | | | |
| | Electron Channel | | Muon Channel | | Electron Channel | | Muon Channel | |
| | 2 jets | 3 jets | 2 jets | 3 jets | 2 jets | 3 jets | 2 jets | 3 jets |
| Signals | | | | | | | | |
| $tb$ | 20 | 8.1 | 20 | 9.4 | 158 | 39 | 133 | 34 |
| Background Sum | 14962 | 3586 | 18610 | 5125 | 78502 | 11526 | 72382 | 11192 |
| Background + Signal | 15021 | 3611 | 18672 | 5156 | 78941 | 11642 | 72764 | 11294 |
| Data | 15021 | 3611 | 18672 | 5156 | 78936 | 11641 | 72762 | 11293 |

**Table 5.13** Pre-tagged event yields after selection.

Estimated range of $\rho$ for selected processes:

Top Quark Production at DØ                                    $\rho \in \langle 0.001, 0.003 \rangle$

Higgs boson search at ATLAS, LHC                             $\rho \simeq 10^{-4} - 10^{-6}$

NOvA: muon antineutrinos $\rightarrow$ electron antineutrinos      18 events over three years
(press release, June 4, 2018)

Estimates of
the learning
set size for
k-NN and IINC
methods in
HEP

František Hakl

Data driven
verification

PAC model
Definitions
Main theorem
Proof

k-NN and IINC
k-NN methods
IINC methods

HEP
separation

Conclusion

Conclusion

- method of learn data size estimation is suggested for very rare processes
- upper bound of the the Vapnik-Chervonenkis dimension for consistent nearest neighbor and IINC like methods is derived
- set size estimation is applicable in applications in which the ratio $\rho$ of positive examples is extremely small

$$\text{NN:} \quad \rho \lessapprox \frac{const. \cdot dim(\bar{X}) \cdot \epsilon}{\log(\# \text{ of pos. examples})}$$

$$\text{IINC:} \quad \rho \lessapprox const. \cdot \epsilon$$