# 048843

Itay Hubara and Zohar Rimon

April 2020

## 1  Testing a simple active algorithm

### 1.1  Synthetic dataset

**a**   Since $\Sigma_+ = \Sigma_- = \Sigma$ and $P(w_-) = P(w_+) = \frac{1}{2}$ we get that the Bayes Classifier is:

$$\frac{1}{2\pi^{k/2} \cdot |\Sigma|^{1/2}} e^{(x-\mu_-)^T \Sigma^{-1}(x-\mu_-)} = \frac{1}{2\pi^{k/2} \cdot |\Sigma|^{1/2}} e^{(x-\mu_+)^T \Sigma^{-1}(x-\mu_+)} \quad (1)$$

Multiplying by $2\pi^{k/2} \cdot |\Sigma|^{1/2}$ and taking log from both sides results with:

$$(x - \mu_-)^T \Sigma^{-1}(x - \mu_-) = (x - \mu_+)^T \Sigma^{-1}(x - \mu_+)$$

Straightforward derivation leads to:

$$2x^T \Sigma^{-1}(\mu_- - \mu_+) = \mu_-^T \Sigma^{-1}\mu_- - \mu_+^T \Sigma - 1\mu_+$$

Thus we get

$$W = \Sigma^{-1}(\mu_- - \mu_+) \; ; \; b = -\frac{1}{2}(\mu_-^T \Sigma^{-1}\mu_- - \mu_+^T \Sigma^{-1}\mu_+)$$

**b**   Here we apply same derivation but for the general case in which $\Sigma_- \neq \Sigma_+$, this results with:

$$\frac{1}{2\pi^{k/2} \cdot |\Sigma_-|^{1/2}} e^{(x-\mu_-)^T \Sigma_-^{-1}(x-\mu_-)} = \frac{1}{2\pi^{k/2} \cdot |\Sigma_+|^{1/2}} e^{(x-\mu_+)^T \Sigma_+^{-1}(x-\mu_+)} \quad (2)$$

$$\frac{|\Sigma_+|^{1/2}}{|\Sigma_-|^{1/2}} e^{(x-\mu_-)^T \Sigma_-^{-1}(x-\mu_-)} = e^{(x-\mu_+)^T \Sigma_+^{-1}(x-\mu_+)} \quad (3)$$

$$\frac{1}{2}log(|\Sigma_+|) - \frac{1}{2}log(|\Sigma_-|) + (x-\mu_-)^T \Sigma_-^{-1}(x-\mu_-) = (x-\mu_+)^T \Sigma_+^{-1}(x-\mu_+) \quad (4)$$

Defining $c1 = \frac{1}{2}log(|\Sigma_+|) - \frac{1}{2}log(|\Sigma_-|)$ and simplifying the equation leads to

$$c1 + x^T(\Sigma_-^{-1} - \Sigma_+^{-1})x + 2x^T(\Sigma_+\mu_+ - \Sigma_-\mu_-) = \mu_+\Sigma_+\mu_+ - \mu_-\Sigma_-\mu_- \quad (5)$$

Thus we results with quadratic classifier: $x^T Ax + x^T b + c$ where $c = c1 + \mu_-\Sigma_-\mu_- - \mu_+\Sigma_+\mu_+$, $A = \Sigma_-^{-1} - \Sigma_+^{-1}$ and $b = 2(\Sigma_+\mu_+ - \Sigma_-\mu_-)$

**c** As can be seen in table 1 ,the first experiment in which the mean is normalized by $sqrt(d)$ the accuracy results is roughly constant while for the second experiment for which $\mu_{+/-} = \pm(1, 1, .., 1)$ we get perfect classification for $d \geq 20$

| Experiment | d=1 | d=20 | d=50 | d=200 |
|---|---|---|---|---|
| $\mu_{\pm} = \pm(1, 1, ...1)/\sqrt{(d)}$ | 84% | 84.5% | 83.5% | 84.6% |
| $\mu_{\pm} = \pm(1, 1, ...1)$ | 84% | 100% | 100% | 100% |

In both of the cases we get the same Bayes estimator (as seen in the first section, with this section's parameters):

$$W = (\mu_- - \mu_+) \ ; \ b = -\frac{1}{2}(\mu_-^T\mu_- - \mu_+^T\mu_+) = 0$$

*We get different Ws but the result is a same classifier.

We can notice that the classifier, as expected is symmetric and is through the origin. more specifically, it does not change with the dimension.

For the normalized case: The distance between the means of the two distributions and the origin stays the same as the dimension increases

$$dist = \sqrt{\sum_{i=1}^{d}(\mu_+ - 0)^2} = \sqrt{\sum_{i=1}^{d}1/d} = 1$$

For the non-normalized case: The distance between the means of the two distributions and the origin grows as the dimension increases

$$dist = \sqrt{\sum_{i=1}^{d}(\mu_+ - 0)^2} = \sqrt{\sum_{i=1}^{d}1} = \sqrt{d}$$

For the normalized case, although the distance of the Gaussian means from the origin is constant and equal to the variance, thus we get constant accuracy of $85\%$ (which makes as for each as for each Gaussian you have 50% that bigger/smaller then the mean and an additional 34% which still on the correct side of the hyper-plane) . For the non-normalized case the distance of the Gaussian's mean from the origin scales with $\sqrt{(d)}$ thus for $d == 20$ we get that the distance from the origin is between $4\sigma - 5\sigma$ thus the classification is perfect.

## 1.2 Test Set Results

*Setting*

Throughout the experiments we used sklearn package which includes a method to optimized the parameters, $SearchGridCV$. This method run hyper parameters tuning using a simple exhaustive grid search (nested for loops) and picks the best parameters based on the cross validation score. The results were averaged over 30 trials (runs in parallel) where in each round the pool/test/val dataset are randomly chosen out of 600

samples. We note that we used the scaled breast cancer dataset no further reprocessing was done to the datasets.

As can be seen in Figure 1. SIMPLE algorithm generalize better on the dataset as it improved its hyper-plane based on the points close to the margin.
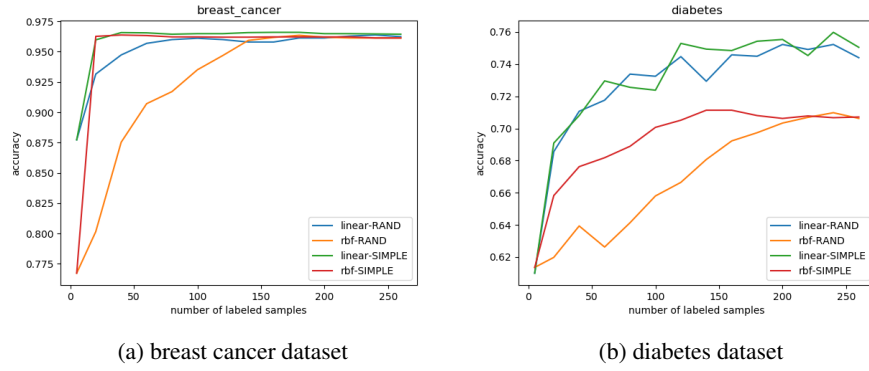


(a) breast cancer dataset  (b) diabetes dataset

Figure 1: Accuracy (over the test set) per number of labeled samples for the "real" datasets

## 1.3   LOO Results

To reduce running time we used cross validation with five segments instead LOO. Here, by using SIMPLE algorithm we pick the hardest points to classify (the one closes to the margin, and test our performance on them. Thus, it falsely appears like that the SIMPLE algorithm perform worse then the RAND. However if we would test the SIMPLE algorithm on the subset of the training set chosen by the RAND algorithm we would see that it is actually better.
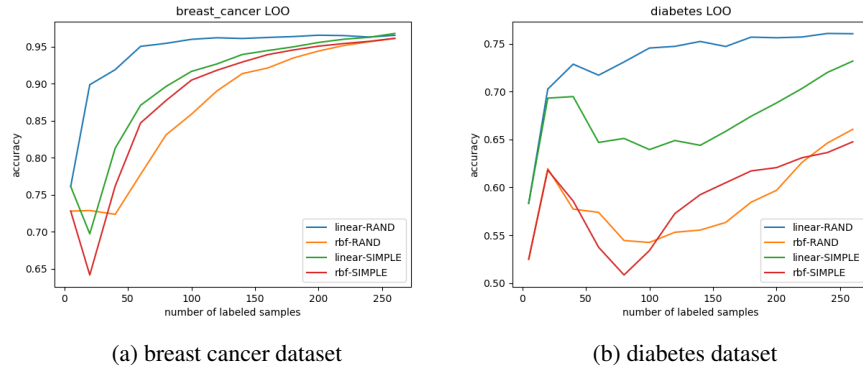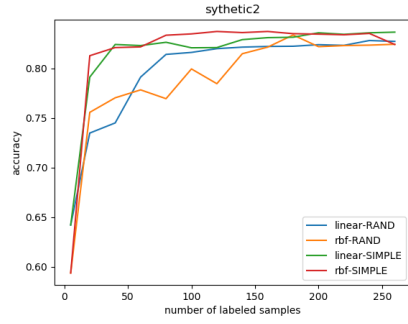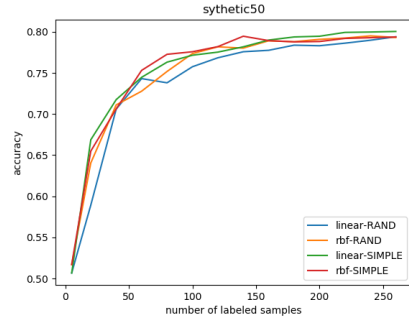
(a) breast cancer dataset          (b) diabetes dataset

Figure 2: Accuracy (using LOO on the labeled samples) per number of labeled samples for the "real" datasets
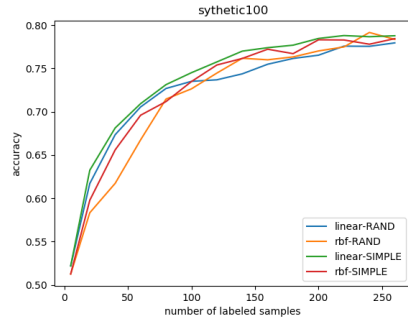
## 1.4 Synthetic dataset results

Here we see the same trends as in the "real" dataset. The SIMPLE algorithm generalize better on the test dataset and if you test the algorithms on the labeled points, since SIMPLE choose the ones closest to the margin its results are inferior, In addition we observe that the accuracy decreases with the problem dimension suggesting that the estimation error increases with the dimension. This means that ,although the Bayes estimator is a linear classifier thus the approximation error is zero the SVM can't find it. Since we know that the samples are not noisy, we simply need to increase the numbers of samples to achieve Bayes error.
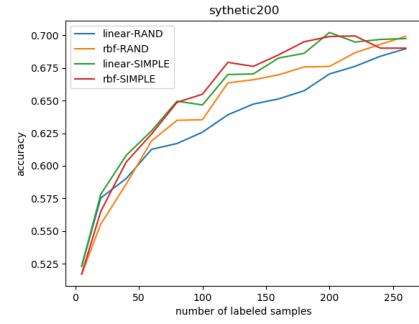
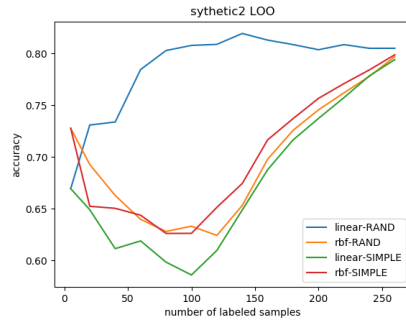(a) Synthetic dataset with d=2

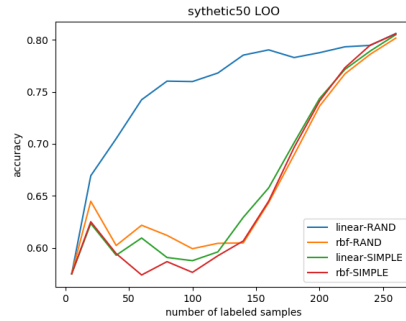(b) Synthetic dataset with d=50

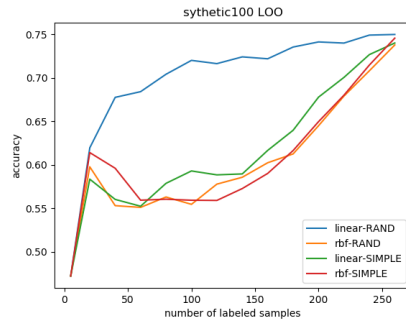(c) Synthetic dataset with d=100

(d) Synthetic dataset with d=200

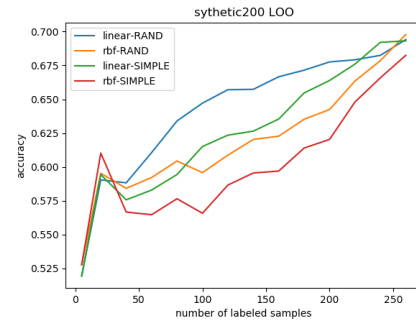Figure 3: plots of synthetic dataset accuracy over the test set v.s number of labels sampled for d=2,50,100,200

(a) Synthetic dataset with d=2



(b) Synthetic dataset with d=50



(c) Synthetic dataset with d=100



(d) Synthetic dataset with d=200

Figure 4: plots of synthetic dataset accuracy using LOO over the labels samples v.s number of labels sampled for d=2,50,100,200

## 1.5 SIMPLE Algorithm Failures

For the linear classifier we'll use the 1D example from the lecture notes in which 95% of the data leads to a hyper-plane on the origin. Since, with high priority we would sample from those examples we wont see the rest of the 5% and if would take a lot of samples until we'll get 97.5 accuracy. For the RBF kernel the example is similar only in 2D with circles. Thus we know that the classifier can reach 97.5% yet it requires almost all the samples to get to that accuracy

(a) breast cancer dataset
(b) diabetes dataset

Figure 5: Simple failure on RBF and Linear classifier

I) נגדיר אב

(1) נראה כי $2k$ נקודות זו הסדר הישר:



```
+----•----•---------------------------•------->
     1    2                          2k
```

... (Hebrew text, handwritten, largely illegible) ...

$$\Rightarrow \boxed{V_u = 2k}$$

(2) נראה אולי $V = 3$ . $x_1 < x_2 < x_3, x_4$
$x_1 < x_2 < x_3 < x_4$

... (Hebrew text) ...

$$\Rightarrow x_1 \in [\alpha, \alpha+1], \quad x_2 \in (\alpha+1, \alpha+2), \quad x_3 \in [\alpha+2, \infty)$$

$$\Rightarrow x_u \in [\alpha+2, \infty)$$

$$\Rightarrow h(x_u) = \ddot{+} \neq \ddot{-}$$

$$V_H < 4$$

נתבונן ב $e$ $E=H$ (כיצד קבוצה) $2 \in$ נקודות 3 $S$ קבוצה

ש(1)shattered!

$$x_1 = \frac{2}{8}, \quad x_2 = \frac{9}{8}, \quad x_3 = \frac{14}{8}$$

נראה כעת כי לכל סוג של תיוג ניתן יהיה למצוא מסווג כזה!

$+ \; + \; + \quad \Rightarrow \quad \alpha = -5 \Rightarrow [-5, -4] \cup [-3, \infty)$

$- \; - \; - \quad \Rightarrow \quad \alpha = 5 \Rightarrow [5, 6] \cup [7, \infty)$

$+ \; + \; - \quad \Rightarrow \quad \alpha = \frac{2}{8} \Rightarrow [\frac{2}{8}, \frac{10}{8}] \cup [\frac{18}{8}, \infty)$

$+ \; - \; + \quad \Rightarrow \quad \alpha = -\frac{3}{8} \Rightarrow [-\frac{3}{8}, \frac{5}{8}] \cup [\frac{13}{8}, \infty)$

$- \; + \; + \quad \Rightarrow \quad \alpha = -1 \Rightarrow [-\frac{7}{8}, \frac{1}{8}] \cup [\frac{9}{8}, \infty)$

$- \; - \; + \quad \Rightarrow \quad \alpha = \frac{13}{16} \Rightarrow [-\frac{13}{16}, \frac{3}{16}] \cup [\frac{19}{16}, \infty)$

$- \; + \; - \quad \Rightarrow \quad \alpha = \frac{1}{2} \Rightarrow [\frac{1}{2}, \frac{3}{2}] \cup [\frac{5}{2}, \infty)$

$+ \; - \; - \quad \Rightarrow \quad \alpha = 0 \Rightarrow [0, 1] \cup [2, \infty)$

$$\boxed{V_H = 3}$$

$\bigcirc =$

$\Rightarrow$ "אם נסווג נקודות ונצייר את השולים" השגיאה החלוט...

ומכאן, כלומר co$^{0}$ אם 5 נקודות אז אבל לא נקודות ל... השגיאה

ובכן $V_H = 5$.

Scanned with CamScanner

טענה 1: ישנה נקודה כזו בגלל תורה הקואורדי:

· שינינו סכום את הדינמיקה של הבוֹ רה כוֹ (A) וזה התנאות כי (B)
גזירת האלוקה. כיצד היינ מופיע כל פא היינה או הדינמיקה?

טענה II: אין נקודה פנימית מ... ...

התכנסות! כל מהיהתכנס של הבוּ רה בתוֹ מתוֹ השלם וזה כן כל קווים
המתכנסות בענף בתוֹ השלם ולכן כל הבו רה התעכב
בתוֹ השלם או ניתן שלא נקודה תהיה נתוֹ...שלי.

טענה II: אין נקודה פנימית ! $P_0$ $P_0$

נעבון| נעקבה השעותוֹ היער וחתוֹ ביקר בין 5 הדינמיקים
(כבר עתוֹ' הן שעונף ותחוֹת)

יתן נקח שפ"ב ב + ולסכן גסוֹ השוֹל בתיך מתֹר

מתוֹ $P_0$ ו... ... ...שים $2P_0$ . פוֹ נוֹ

הדינמיק שנתֹ זוֹנ לשווות שתתברך ... שפ שפ...

נסבר יך 3. הדינמיק שנתֹ כ:/$P_1$,/$P_2$/$P_3$,/ בסֹ...
זיע יך! $(P_3) \times (P_2) \times (P_4) \times (P_3)$ (כֹבוֹ עתֹ ... אין ...)

...לוֹ ...יין וֹין נקוֹ בתוֹ הבוֹרה הקואנטוֹ!

היֹר התֹר בין / $P_1$ לכן $P_1$ בֹיק כֹבוֹ מתֹ 2 $P_0$. (A)

עתוֹ עוֹ : $-P_3'$,$+P_2'$, $P_0':-P_1'$ בֹיק שֹיֹ יֹתֹר מֹתֹ

... 1 כֹ $P_0$ ושֹו וֹקֹ כֹ שֹ"ינֹ (A).

⟹בֹ נתֹ כֹ ... סֹיֹ ... ...

מֹתֹ בֹכֹ... ...קֹוֹם וֹ נֹכֹ... יֹ... השֹוֹ... בֹיֹ... וֹחֹוֹ בֹין ...ֹבֹין

יֹ...יֹתֹ הֹנֹ...ֹ הֹ...ֹ וֹ...ֹן כֹבֹ.

(ii) יֹוֹ ... כֹ נֹ... ... ... ... יֹ... בֹ...ֹ (... ...ֹ בֹ...)

נתיר מהרצאות הנקודות עד היינך ייי ונשק טנו.

שיווין $x(\beta_1)=x(\beta_1)$   ההומתי נשיון נבונ.

תזכה!

$$P(X>t) \le Ce^{-2nt^2}$$

$$E(X^2) = \int_0^\infty P(X^2 \ge t)dt = \int_0^2 P(X^2 \ge t)dt + \int_2^\infty P(X^2 \ge t)$$

$$\underset{P\le 1}{\le} 2 + \int_2^\infty P(X^2 \ge t) = 2 + \int_2^\infty P(X \ge \sqrt{t})dt$$

$$\le 2 + C\int_2^\infty e^{-2nt}dt = 2 + \frac{Ce^{-2n\cdot 2}}{2n}$$

נבחר $2 = \log(C)/n$ נקבל!

$$E[X^2] \le \frac{\log(C)}{2n} + \frac{1}{2n} = \frac{\log(Ce)}{2n}$$

ומכאן! מקבל

$$E[X] \le \sqrt{E[X^2]} \le \sqrt{\frac{\log(Ce)}{2n}}$$

(1)נוכיח כי $d$ על $B$ מטריקה!

$$d(h,h') = P\{h(x) \ne h'(x)\} \ge 0$$

כי הסתברות

$$d(h,h') = d(h',h) = P(h(x) \ne h'(x))$$

אזיי אם $h = h'$ כי

$$d(h,h) = P(h(x) \ne h(x)) = C$$

אזיי אם $h \ne h'$ כי $e$ קיים $x$ שעבורו $h(x) \ne h'(x)$ (בהסתברות חיובית מסויימת)

$$d(h,h') \ne 0$$

$$d(h,h') + d(h',h'') = P(h(x) \neq h'(x)) + P(h'(x) \neq h''(x))$$

$$\geq P(h'(x) \neq h''(x) \wedge h(x) = h'(x))$$
$$+ P(h(x) \neq h''(x) \wedge h(x) \neq h'(x))$$
$$= P(h(x) \neq h''(x)) = d(h,h'')$$

$$OIS(B(h, \varepsilon/\sqrt{d})) = \qquad (2)$$

$$= \{x \in \mathbb{R}^d ; \exists h, h' \in B(h, \varepsilon/\sqrt{d}) \text{ with } h(x) \neq h'(x)\}$$

$$B(h, \varepsilon/\sqrt{d}) = \{h' \in H : d(h,h') \leq \frac{\varepsilon}{\sqrt{d}}\}$$

<span dir="rtl">אולי</span> $x$ <span dir="rtl">אולי</span> $N\leq$ <span dir="rtl">בכל נקודה אחיד מתוך כדור היחידה ב</span> $\mathbb{R}^d$ :

$$d(h,h') = P\{h(x) \neq h'(x)\} = \frac{V(h(x) \neq h'(x) : x \in \text{unit-ball})}{V_d} =$$

<span dir="rtl">גודל</span> $\nearrow k$     <span dir="rtl">נפח כל האזורים</span>    $V_d \searrow$    <span dir="rtl">נפח unit ball</span>

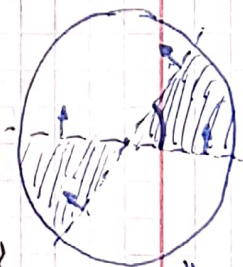$$= \frac{V(\text{sign}(\omega \cdot x) \neq \text{sign}(\omega' \cdot x) : x \in \text{unit-ball})}{V_d} =$$

<span dir="rtl">ראינו בהרצאה שיחס הנפחים</span>   $\dfrac{\text{arccos}(\omega \cdot \omega')}{\pi}$

<span dir="rtl">ניקח בה"כ:</span> $\omega = (1, 0, 0, \cdots)$

$$\Rightarrow B(h, \varepsilon/\sqrt{d}) = \{h' \in H ; \frac{\text{arccos}(\omega \cdot \omega')}{\pi} \leq \frac{\varepsilon}{\sqrt{d}}\} =$$

$$= \{h' \in H ; \frac{\text{arcos}(\omega_1')}{\pi} \leq \frac{\varepsilon}{\sqrt{d}}\}$$

$$= \{h \in H : \omega_1' \geq \cos(\frac{\pi \varepsilon}{\sqrt{d}})\}$$

$$\Rightarrow \boxed{OIS[B(h, \frac{\varepsilon}{\sqrt{d}})] = \{x \in \text{unit-ball} ; |x_1 \omega_1| \leq \sin(\pi \frac{\varepsilon}{\sqrt{d}})\}}$$

<span dir="rtl">המישור החותך</span> $\downarrow \overset{V}{=}$

$$=$$

<span dir="rtl">והמישורים</span>