$$S = (z_1, ..., z_m)$$

$$z_i := (x_i, y_i) \qquad \text{דוגמה בודדת ...}$$

$$z_i' \qquad \text{דוגמה נוספת ...}$$

$$S^{(i)} = (z_1, ..., z_{i-1}, z_i', z_{i+1}, ..., z_m)$$

$$U(m) \sim U\{1, m\}$$

$$f_S(w) = \frac{1}{m} \sum_i^m \ell(w, z_i) + \lambda \|w\|^2$$

$$\hat{w} = \arg\min f_S(w)$$

סימון:

$$L_D(\hat{w}) = E_{(z) \sim D}\, \ell(\hat{w}, z)$$

$$L_S(\hat{w}) = \frac{1}{m} \sum \ell(\hat{w}, z_i)$$

רוצים:

$$E_{S \sim D^m}\left[L_D(\hat{w}) - L_S(\hat{w})\right] = E_{(S, z') \sim D^{m+1},\ i \sim U(m)}\left[\ell(\hat{w}^{(i)}, z_i) - \ell(\hat{w}, z_i)\right]$$

* ...

* ...

C ... $loss$

הוכחה:

$$E_S\{L_D(\hat{w})\} = E_{S,z}\left(\ell(\hat{w}, z_i)\right) = E_{S,z'}\left(\ell(\hat{w}^{(i)}, z_i)\right) \qquad \text{①}$$

$$E_S\{L_S(\hat{w})\} = E_{S,i}\left\{\ell(\hat{w}, z_i)\right\} \qquad \text{②}$$

... ההפרש (1)-(2) ...

$$Y = f(x) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

נזכיר: מודל $\hat{f}(x)$ שנאמד על בסיס הנתונים $\quad E\{(Y - \hat{f}(x))^2\}$

ונזכיר:

$$Bias\{\hat{f}(x)\} = E\{\hat{f}(x) - f(x)\}$$

$$Var\{\hat{f}(x)\} = E[\hat{f}(x)^2] - (E[\hat{f}(x)])^2$$

Bias = נותן לנו מדד עד כמה המודל מסוגל ללכוד את הקשר האמיתי בין
המשתנים. ככל שזה גבוה יותר Bias (מדד גבוה) אזי המודל לוכד פחות טוב
המודל מסוגל 10000 פעמים וכו' מנסה ללכוד את הקשר בין המשתנים
זה ביטוי לחוסר יכולת (under fitting)

Variance = נותן לנו מדד עד כמה המודל מסוגל ללכוד שינויים קטנים בנתון.
variance גבוה זהו מצב שבו המודל לוכד שינויים קטנים וכו' עד כדי כך
זה מצב שבו גם את הרעש. מנסים לומר שהמודל לומד יותר מדי,
variance גבוה מצב של overfitting

נסתכל בגרפים:



| High bias | מצב של מאוזן | low Bias |
| low varians | בין השניים | High variance |

• נשים לב שככל ה Bias גבוה יותר, המודל נעשה יותר פשוט (ליניארי), ואכן
מתקיים שככל שהמודל יותר פשוט כך נקבל bias גבוה יותר. בנוסף
נראה ש x mse נראה גם וכן מתקיים mse גבוה יותר.
• מצד שני $\sigma = $ variance גבוה, יותר הופך להיות מורכב, וכן
וכן זה לשים לב שככל שהמודל יותר מורכב כך יהיה לנו variance
גבוה יותר ומצד שני bias גבוה יותר.

$$E\left[(y-\hat{f}_{(x)})^2\right] = \left(Bias\{\hat{f}_{(x)}\}\right)^2 + Var\left[\hat{f}_{(x)}\right] + \sigma^2 \quad \rightarrow$$

נסמן $f=f_{(x)}$, $\hat{f}=\hat{f}_{(x)}$ .

1) $E(x^2) = Var(x) + (E(x))^2$

. נניח $\hat{f}-\theta$ ביטוי של $Var$ -ל

. נניח ש $f-\theta$ זהו ביטוי $covariance$ זהו ביטוי

$$E\{f\} = f$$

$\Rightarrow$ 2) $y = f + \varepsilon$ , and $E(\varepsilon) = 0 \Rightarrow E(y) = E(f+\varepsilon) = E(f) + E(\varepsilon) = f$

$\uparrow$ של ביטוי השגיאה

$$Var\{\varepsilon\} = \sigma^2$$

. של ביטוים:

נביט בביטוי השמאלי $Var$:

3) $Var\{y\} = E\left[(y - E(y))^2\right] = E\left[(y-f)^2\right] = E\left[(f+\varepsilon-f)\right] = E(\varepsilon^2)$ ,

$\left(1\right)$ לפי $= Var(\varepsilon) + (E(x))^2 = \sigma^2$

$0''$ של ביטוי השגיאה

. נפתח את אגף ימין כך:

$$E\left[(y-\hat{f})^2\right] = E\left[(f+\varepsilon-\hat{f})^2\right] = E\left\{(f+\varepsilon-\hat{f}+E(\hat{f})-E(\hat{f}))^2\right\}$$

$$= E\left[(f-E(\hat{f}))^2\right] + E(\varepsilon^2) + E\left((E(\hat{f})-\hat{f})^2\right) + 2E\left((f-E(\hat{f}))\right)\cdot E(\varepsilon)$$

$0''$

$$+ 2E\left((f-E(\hat{f}))\cdot E((E(\hat{f})-\hat{f}))\right)$$

$0''$

$$+ E(\varepsilon)\cdot E\left((E(\hat{f})-\hat{f})\right)$$

$0$

$$= E\left[(f-E(\hat{f}))^2\right] + E(\varepsilon^2) + E\left((E(\hat{f})-\hat{f})^2\right)$$

$$= (f-E(\hat{f}))^2 + \sigma^2 + Var(\hat{f})$$

$$= Bias(\hat{f}_{(x)})^2 + \sigma^2 + Var(\hat{f}) \qquad \text{מ.ש.ל}$$

**2.** א. ‎הגדרת ‎... ‎כי ‎בכל ‎ב"ן ‎מאיץ ‎בקבוצה ‎‎

‎... ‎שלוין ‎של ‎ŵ, ‎בחוללה ‎... ‎ניתן ‎על ‎אחד ‎

‎הלינין, ‎ו... ‎שלוין ‎של ‎ŵ ‎על ‎גוף ‎... ‎כשר,

‎"ל..." ‎... ‎H ‎כך ‎... ‎... ‎ ‎...,

‎... ‎של ‎... ‎אלגוריתם ‎ההכרעות ‎כיבוי, ‎...,

‎... ‎... ‎(עצ) ‎... ‎וכבוי ‎... ‎גרון/גלינ.

⇐ ‎כל ‎... ‎של, ‎... ‎... ‎... ‎על ‎...

ג.

א. ‎... ‎שלוש ‎...:

(i) ‎ $loss(z,w) = loss(z,u) = 0$

‎אזי: ‎ $|loss(z,w) - loss(z,u)| = 0 \le |z|\,\|w-u\|$

(ii) ‎ $loss(z,w) = 1 - y\langle w,x\rangle$ , $loss(z,u) = 1 - y\langle u,x\rangle$

$|loss(z,w) - loss(z,u)| = |1 - y\langle w,x\rangle - (1 - y\langle u,x\rangle)|$

$= |y|\cdot|\langle w,x\rangle - \langle u,x\rangle| = |(w-u)\cdot x| \le |w-u|\cdot|x|$

(iii) ‎בה"כ ‎ $loss(z,u)$

$loss(z,w) = 1 - y\langle w,x\rangle > 0$

$loss(z,u) = 0 > 1 - y\langle u,x\rangle$

‎ונקבל:

$|loss(z,w) - loss(z,u)|$

$= |1 - y\langle w,x\rangle - 0| \le |1 - y\langle w,x\rangle - (1 - y\langle u,x\rangle)| \le |w-u|\cdot|x|$

‎ ‎$\underbrace{\phantom{1-y\langle w,x\rangle}}_{0}$ ‎$\uparrow$ ‎$\uparrow$ ‎$\nearrow$

‎לפי ‎סעיף (ii)

ב. ‎ניתן ‎להפוך ‎כי ‎כאשר ‎כל ‎... ‎x ‎... ‎... ‎ל"ל.

‎ה-domain ‎התפוס ‎צ', ‎ושל ‎כלל ‎... ‎C ‎, soft-SVM

‎גודל ‎C-... ‎... ‎... ‎... ‎... ‎.

ב. ג. נגדיר $loss(z,w)=log(1+e^{-y\langle w,x\rangle})$

נגדיר ב- $f(t)=log(1+e^{-t})$

כיוון ש- $e^{-t}>0$ לכל $t\in\mathbb{R}$, $f(t)$ גזירה ב- $\mathbb{R}$

ולכן, על פי משפט ערך הביניים (הממוצע) קיים $c$ כך שמתקיים

$$f'(c)=\frac{f(a)-f(b)}{a-b}\qquad \forall a,b\in\mathbb{R}$$
$$\text{כאשר } c \text{ כך ש- } b<c<a$$

$$f'(t)=\frac{1}{1+e^{t}}\cdot(-e^{-b})=-\frac{e^{-b}}{1+e^{-b}}$$

סכום הנגזרות (או סכום חיובי) , לכן צריך לבדוק ולהראות

$$|f'(b)|<1\qquad\forall t\in\mathbb{R}$$

ומתקיים

$$\text{I}\quad f(a)-f(b)\le |a-b|$$

$$(\text{1-ליפשיץ})$$

נראה שכעת הפונקציה $loss(z,w)$ היא ליפשיץ בכבוד , $x,w\in\mathbb{R}$

נגדיר $f\circ g(z,w)$ , נסמן $g(z,w)=yxw$

לכן על פי כלל השרשרת נקבל

$$\frac{\partial\, logloss}{\partial w}=\frac{\partial f}{\partial t}\frac{\partial g}{\partial w}\le yx$$

ובכך נקבל ליפשיץ ל- $c$.

$$logloss(z,w)-logloss(z,u)\le |w-u||yx|=|w-u||x|$$

ד. נוכיח ליפשיץ ב- נבחר $x$ כך שיהי אורתוגונלי ל- $w$ ושהם
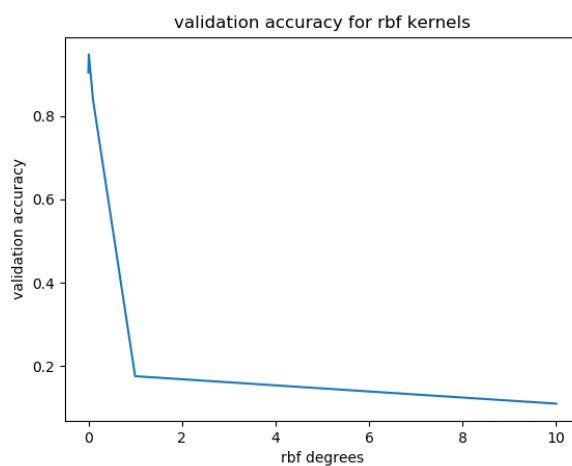ומעלים כך שיהי אורתוגונלי אזי $\emptyset$ אינו ריקה ליפשיץ לכל הכיוון

4)

Results:

linear: [1.0, 0.90390625, 0.913125]
poly_2: [1.0, 0.9534374999999999, 0.960625]
poly_3: [1.0, 0.94625, 0.950625]
poly_4: [1.0, 0.9262499999999999, 0.94125]
poly_5: [1.0, 0.9018750000000001, 0.920625]
poly_6: [1.0, 0.8728125, 0.895625]
poly_7: [1.0, 0.8417187500000001, 0.86875]
poly_8: [1.0, 0.813125, 0.845625]
poly_9: [1.0, 0.7834375, 0.8125]
rbf_0.001: [0.9140625, 0.9040625, 0.903125]
rbf_0.01: [0.9776953124999999, 0.9475000000000001, 0.951875]
rbf_0.1: [1.0, 0.83890625, 0.85625]
rbf_1: [1.0, 0.17625000000000002, 0.165625]
rbf_10: [1.0, 0.11046875, 0.105]

the poly_2 model was best on cross validation with error 0.9534374999999999

the poly_2 model was best on test data with error 0.960625

$z_i = (x_i, y_i)$ , $S = (z_1, ... z_m)$ , $z'$ — נקודה חדשה

$S^{(i)} = (z_1, ... z'_i ... z_m)$ (דוגמה):

* נתבונן במשל $L$ אשר $r$ הנחה. נבחן $r$ כלשהו אשר משנה נקודה אחת $z_i$

* input ל־$L$ אשר $G$ מחזיר מודל $r$. loss. $S^{(i)}$ שבו החלפנו $z_i$

* במ... נקבה $r$ מודל $z'$ במקום $S - S$

$$\hat{w} = \text{argmin}\, f_S(w) \quad \text{כאשר}\quad f_S(w) = L_S(w) + \lambda \|w\|^2 \quad \text{(אימון)}$$

$\hat{w}$ — נראה כי $f_S$ הוא ... עבור $\hat{w} = $ ... כלומר ...

... $S$ ...

$$f_S(v) - f_S(w) = L_S(v) + \lambda |v|^2 - \left(L_S(w) + \lambda |w|^2\right) \quad \text{נשתמש}\ v, u\ ...$$

$$= L_{S^{(i)}}(v) + \lambda \|v\|^2 - \left(L_{S^{(i)}}(u) + \lambda \|u\|^2\right) + \frac{\ell(v, z_i) - \ell(u, z_i)}{m} + \frac{\ell(u, z'_i) - \ell(v, z'_i)}{m}$$

— ... החלפנו $z'$ במקום $z_i$ ...

$$L_S(w) = \frac{1}{m} \sum_{i=1}^{m} \ell(w, z) = L_{S^{(i)}} + \frac{\ell(w, z_i)}{m} - \frac{\ell(u, z'_i)}{m}$$

— נכון לכל ... $u, v$ ...

. נבחר ... $u = \hat{w}$ , $v = \hat{w}^{(i)}$ ...

$$L_{S^{(i)}} w + \lambda \|v\|^2 \qquad 0-v$$

...

$$f_S(\hat{w}^{(i)}) - f_S(\hat{w}) \leq \frac{\ell(\hat{w}^{(i)}, z_i) - \ell(\hat{w}, z_i)}{m} + \frac{\ell(\hat{w}, z'_i) - \ell(\hat{w}^{(i)}, z'_i)}{m}$$

$$\eta\,\|\hat{W}^{(i)} - \hat{W}\|^2 \leq \frac{\ell(\hat{W}^{(i)}, z_i) - \ell(\hat{W}, z_i)}{m} + \frac{\ell(\hat{W}, z_i') - \ell(\hat{W}^{(i)}, z_i')}{m}$$

וזה נכון עבור כל $i$. כעת נסתכל על כל אחד מהביטויים. נשים לב כי [...] על הסכום של הערכים של [...] ואז גם על כל [...] לא נסתכל על הסכום [...] (סכום על $z_i$ ועוד סכום על $z_i'$). כעת נעשה [...] כדי לקבל overfit על קבוצת האימון.

1. נסתכל על ה-loss על $\ell_d$, נדע כי נקבל:

*)  $\ell(\hat{W}^{(i)}, z_i) - \ell(\hat{W}, z_i) \leq \rho\,|\hat{W}^{(i)} - \hat{W}|$

ובאותו אופן:

2)  $\ell(\hat{W}, z_i') - \ell(\hat{W}^{(i)}, z_i') \leq \rho\,|\hat{V}^{(i)} - \hat{W}|$

נציב את הביטוי $1+2$ לביטוי למעלה ונקבל:

$$\eta\,\|\hat{W}^{(i)} - \hat{W}\|^2 \leq \frac{\rho\,|\hat{W}^{(i)} - \hat{W}|}{m} + \frac{\rho\,|\hat{W}^{(i)} - \hat{W}|}{m} = \frac{2\rho\,|\hat{W}^{(i)} - \hat{W}|}{m}$$

$$\Rightarrow \|\hat{W}^{(i)} - \hat{W}\| \leq \frac{2\rho}{\eta m} \quad\Rightarrow\quad \ell(\hat{W}^{(i)}, z_i) - \ell(\hat{W}, z_i) \leq \frac{2\rho}{\eta m}$$

=) אם אנחנו מחפשים סכום [...] ביטוי זה [...] ה-loss [...] מראה [...] עם הערך של $\rho$. כעת נוכל [...] שלכל אלגוריתם [...] את ההפסד על [...] זה [...]. אם [...] נקבל כי [...], אז נראה שבכל [...] אשר שומרים עליהם לקבל ערך [...] ואם נראה אלגו' כזה.