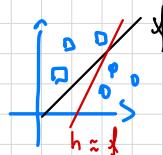




SUMMARY



PAC LEARNING :

LET f BE THE CORRECT CLASSIFIER ON

CORRECT ONLY FOR TRAINING SET

\rightarrow WE SHOULD FIND A FUNCTION $h \approx f$

DATASET

$$X = \{x_1, \dots, x_n\}$$

CORRECT FOR ALL DATASET

→ AN APPROXIMATION

ONE POSSIBLE DEFINITION IS THE ERROR OF h WITH REFERENCE TO f :

x_i DRAWN WITH

$P = P_x$:

$$L_{D, f}(h) := P_{x \sim D} [h(x) \neq f(x)] := D(\{x \in X : h(x) \neq f(x)\})$$



D : A DISTRIBUTION OVER X \rightarrow GIVEN $A \subset X$, $D(A) = P[x \in A]$

SIMPLE DATA GENERATION MODEL:

i.i.d.: $\forall x_i$ is sampled INDEPENDENTLY according to D

REALIZABILITY: $\forall i \in [m], y_i = f(x_i)$

WE CANNOT HOPE TO FIND $h / L_{D, f}(h) = 0$

PROOF:

$\forall \varepsilon \in (0, 1)$, TAKE $X = \{x_1, x_2\}$, $D: D(\{x_1\}) = 1 - \varepsilon$, $D(\{x_2\}) = \varepsilon$

THE PROBABILITY OF NOT SEEING x_2 AMONG m EXAMPLES

$$\rightarrow (1 - \varepsilon)^m \approx e^{-\varepsilon m} \quad \text{→ PROPERTY}$$

\rightarrow IF $\varepsilon \cdot m < 1$; WE ARE UNLIKELY NOT TO SEE x_2 AT ALL

\hookrightarrow RELAXATION $n^{\circ} 1$: WE'D BE HAPPY TO FIND $L_{D, f}(h) \leq \varepsilon$ → USER SPECIFIED

N.B.: $L_{D, f}(h) = 0 \rightarrow h = f$, \because IS TRUE ONLY ON THE PROBABILITY SET CONSIDERED FOR TRAINING, NOT ALL DATA

• NO ALGORITHM CAN GUARANTEE $L_{D,f}(h) \leq \varepsilon$

↳ RELAXATION n°2 . WE ALLOW ALGORITHM TO FAIL $\leq \varepsilon$

WITH PROBABILITY $S \in (0,1)$ \rightsquigarrow LIKE C.I.

$$1-S \leq \varepsilon$$

• PAC LEARNING := PROBABLY APPROXIMATE CORRECT LEARNING

- THE LEARNER DO NOT KNOW D AND \neq
- THE LEARNER RECEIVES ε (ACCURACY PARAMETER) AND S (CONFIDENCE PARAMETER)
- THE LEARNER CAN ASK FOR TRAINING DATA S, WITH m SAMPLES / m(.) ε, S

- THE LEARNER SHOULD OUTPUT AN HYPOTHESIS h WITH PROBABILITY OF AT LEAST $1-S$ / IT HOLDS THAT $L_{D,f}(h) \leq \varepsilon$

$$\rightarrow P_{S \sim D^m} [L_{D,f}(h) \leq \varepsilon] \geq 1-S$$

$S \downarrow \rightarrow P$ OF SUCCESS

$1-S \rightarrow P$ OF FAILURE

• THEOREM (NO FREE LUNCH):

Fix $S \in (0,1)$, $\varepsilon < 1/2$:

$\rightarrow \forall$ LEARNER A AND TRAINING SIZE m:

$\rightarrow \exists D,f$ / WITH $p \geq S$, OVER TRAINING DATA S OF m EXAMPLES
IT HOLDS THAT $L_{D,f}(A(S)) \geq \varepsilon$

• EXPLANATION:

• \forall LEARNER $\rightarrow \exists$ A FUNCTION WHICH DO NOT SUCCEED IN $L_{D,f} \leq \varepsilon$

• \nexists A UNIQUE LEARNER : EVERY TIME A SPECIFIC ONE IS NEEDED

• THE AVG RESULT OF ALL CLASSIFIERS IS WORST THAN RANDOM GUESSING

$$\rightarrow L_{D,f}(\text{RANDOM GUESSING}) = 1/2$$

• IF NO KNOWLEDGE IS INSERTED $\rightarrow \nexists$ SOMETHING BETTER THAN RANDOM GUESSING

• PRIOR KNOWLEDGE :

(NO FREE LUNCH THEOREM)

IF LEARNER KNOWS $H \subset \mathcal{Y}^X$ \rightarrow PAC LEARNING IS POSSIBLE IF $H \neq \mathcal{Y}^X$
 \hookrightarrow HYPOTHESIS CLASS OF \mathcal{Y}

• PAC LEARNING OF FINITE HYPOTHESIS CLASS :

ASSUME H : FINITE HYPOTHESIS CLASS (e.g. all $h: X \mapsto \mathcal{Y}$)

- EMPIRICAL RISK MINIMIZATION (ERM) : $\rightarrow h = \text{ERM}_H(S)$

ERM LEARNING RULE $\left\{ \begin{array}{l} \cdot \text{INPUT: TRAINING SET } S = (x_1, r_1), \dots, (x_m, r_m) \\ \cdot \text{EMPIRICAL RISK: } L_S(h) = \frac{1}{m} |\{i : h(x_i) \neq r_i\}| \quad \{ \text{0-1 loss} \\ \cdot \text{OUTPUT: ANY } h \in H / h \text{ MINIMIZES } L_S(h) \end{array} \right.$

• THEOREM:

Fix ϵ, δ

$$\text{IF } m \geq \frac{\log(|H|/\delta)}{\epsilon^2} \rightarrow \forall D, f, \text{ with } p \geq 1-\delta; L_{D,f}(\text{ERM}_H(S)) \leq \epsilon$$

(PROOF NOT DONE FOR NOW)

OVER TIME CHOICE
OF S SIZE m

$h / \text{SATISFIES ERM}$

$$\Rightarrow m \geq \frac{\log(|H|/\delta)}{\epsilon^2} \geq m_H(\epsilon, \delta)$$

• PAC LEARNABILITY :

AN HYPOTHESIS CLASS H IS PAC LEARNABLE

SAMPLE COMPLEXITY OF LEARNING H

IF \exists FUNCTION $m_H: (0, 1)^2 \mapsto \mathbb{N}$ AND A LEARNING ALGORITHM /:
 \hookrightarrow FOR ϵ, δ

$$\cdot \forall \epsilon, \delta \in (0, 1)$$

$$\cdot \forall D \text{ over } X, \text{ AND } \forall \text{ LABELING FUNCTION } f: X \mapsto \{0, 1\}$$

WHEN RUNNING THE ALGORITHM ON $m \geq m_H(\epsilon, \delta)$ iid EXAMPLES

GENERATED BY D AND LABELED BY f \hookrightarrow SAMPLE COMPLEXITY: min # INPUT DATA TO ASSURE

\rightarrow ALGORITHM RETURNS $h / \text{WHICH } p \geq 1-\delta: L_{D,f}(h) \leq \epsilon$

• COROLLARY: LET H BE A FINITE HYPOTHESIS CLASS:

$\rightarrow H$ IS PAC LEARNABLE WITH $m_H(\epsilon, \delta) \leq \frac{\log(|H|/\delta)}{\epsilon^2}$ / m_H OBTAINED
 USING ERM_H LEARNING RULE

• VC DIMENSION :

• SET $C = \{x_1, \dots, x_{|C|}\} \subset X$ SET OF POINTS, WITH CARDINALITY $|C|$

• SET H_C BE A RESTRICTION OF H TO SET C : $H_C = \{h_c \cdot h \in H\}$

/ $h_c : C \mapsto \{0, 1\}$ AND $h_c(x_i) = h(x_i), \forall x_i \in C$
 ↳ IT'S A FUNCTION THAT MAP THE EACH POINT $\in C$ TO $\{0, 1\}$

→ h_c CAN BE REPRESENTED AS A VECTOR : $(h(x_1), \dots, h(x_{|C|})) = \{\pm 1\}^{|C|}$

$$\rightarrow |H_C| \leq 2^{|C|}$$

• WE SAY THAT H "SHATTERS" C IF $|H_C| = 2^{|C|}$

↳ IF A SET OF FUNCTIONS $\in H_C$ IS ABLE TO PRODUCE EVERY COMBINATION OF $\{\pm 1\}^{|C|}$ OUTPUTS → H SHATTERS C

• $VC_{dim}(H) = d = \sup \{ |C| : H \text{ shatters } C \}$

↳ IT IS THE HIGHEST CARDINALITY OF A SET OF C / H SHATTERS THE SET C

• IN ORDER TO SHOW THAT $VC_{dim}(H) = d$:

1. SHOW THAT \exists A SET C OF SIZE $|C| = d$ / C IS SHATTERED BY H → $VC_{dim}(H) \geq d$
 ↳ JUST 1 SET IS NEEDED

2. \forall SET $C / |C| = d+1$, C IS NOT SHATTERED BY H → $VC_{dim}(H) < d+1$

$$\rightarrow ① + ② : VC_{dim}(H) = d$$

ex.

$$H = \{x \mapsto \text{sign}(x - \theta) : \theta \in \mathbb{R}\}, \quad X \in \mathbb{R}$$

→

• SHOW THAT $\{0\}$ IS SHATTERED :

$$C = \{0\} \rightarrow |C| = 1$$

$$\rightarrow 2^{|C|} = 2^1 = 2$$

$$\rightarrow H_C = \{\text{sign}(-\theta) : \theta \in \mathbb{R}\}, \quad x=0 \mapsto \{\pm 1\} \rightarrow H \text{ shatters } C$$

$$\rightarrow VC_{dim}(H) \geq 1$$

- SHOW THAT ANY 2 POINTS CANNOT BE SHATTERED :



$$\rightarrow \text{VC dim}(H) \leq 2 \rightarrow \text{VC dim}(H) = 1$$

20.

$$H = \{h_{a,b} : a, b \in \mathbb{R}\}, h_{a,b}(x) = 1 \Leftrightarrow x \in [a, b], x \in \mathbb{R}$$

- SHOW THAT $\{0, 1\}$ IS SHATTERED :

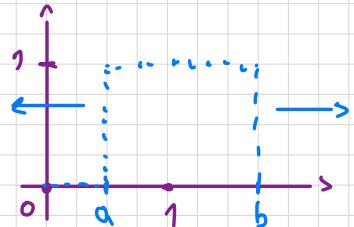
$$\cdot \text{show that } \{0, 1\} \mapsto \{ \pm 1 \}^{|C|=2} \rightarrow \begin{matrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{matrix}$$

$$\cdot (0, 0) : \text{choose } a > 1 \rightarrow \text{as. } h_{2,3}$$

$$\cdot (0, 1) : \text{choose } a \in [0, 1], b > 1 \rightarrow \text{as. } h_{0,5,3}$$

$$\cdot (1, 0) : \text{choose } a < 0, b \in [0, 1] \rightarrow \text{as. } h_{-1,0,5}$$

$$\cdot (1, 1) : \text{choose } a < 0, b > 1 \rightarrow \text{as. } h_{-1,2}$$



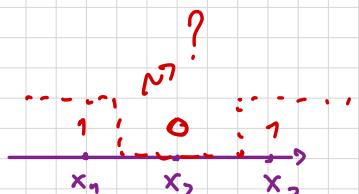
- SHOW THAT ANY 3 POINTS CANNOT BE SHATTERED :

$$\rightarrow \text{THE TRIPLET } (1, 0, 1) \text{ IS NOT POSSIBLE}$$

$$\rightarrow \text{VC dim}(H) \geq 2, \text{VC dim}(H) \leq 3 \rightarrow \text{VC dim} = 2$$

as. VC dim OF MGF SPACES

(SEE LESSONS NOTES, PAGE 79)



• THE FUNDAMENTAL THEOREM OF PAC LEARNING:

LET :

- H BE A FINITE HYPOTHESIS CLASS OF FUNCTIONS $X \mapsto \{0, 1\}$

- LOSS FUNCTION : 0-1 LOSS

→ TFAE: ~ 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 A SUCCESSFUL PAC LEARNER FOR H

6. H HAS A FINITE VC-DIM = $d \leq \infty$

- $\rightarrow \text{VC dim}(H) = d \rightarrow H \text{ IS PAC LEARNABLE}$

• THEOREM:

LET :

- H BE A FINITE HYPOTHESIS CLASS OF FUNCTIONS $X \mapsto \{0, 1\}$

- LOSS FUNCTION : 0-1 LOSS

- VC DIM = $d \leq \infty$

$\rightarrow \exists C_1, C_2$ CONSTANTS, SUCH THAT:

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

AND $m_H(\epsilon, \delta)$ IS ACHIEVED BY THE ERM LEARNING RULE

(NO PROOF)

• LEMMA:

LET H BE AN HYPOTHESIS CLASS / $\text{VC dim}(H) = d < \infty$

$$\rightarrow \forall C \subset X / |C| = m > d+1 \rightarrow |H_C| \leq \left(\frac{e^m}{d}\right)^d \approx K \cdot m^d$$

GENERAL LEARNING MODEL AND BIAS COMPLEXITY TRADEOFF:

THE GENERAL PAC MODEL:

RELAXING THE REALIZABILITY ASSUMPTION - AGNOSTIC PAC LEARNING:

SO FAR WE ASSUMED THAT Y GENERATED BY $f \in \mathcal{H}$

→ THIS ASSUMPTION IS TOO STRONG!

"TARGET GENERATING" FUNCTION $\xrightarrow{\text{ }}$ "DATA-GENERATING DISTRIBUTION"

↳ AGNOSTIC PAC LEARNING:

• D over $X \times Y$

• RISK: $L_D(h) \stackrel{\text{def}}{=} P_{(x,r) \sim D} [h(x) \neq r] \stackrel{\text{def}}{=} D(\{(x, r) : h(x) \neq r\})$

• "APPROXIMATELY CORRECT" NOTION: $L_D(A(S)) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$

↳ A : LEARNER, S : TRAINING SET

BINOMIAL LOSS OF CLASS

LOSS FUNCTION:

LET $Z = X \times Y$. GIVEN $h \in \mathcal{H}$ AND AN EXAMPLE $(\bar{x}, \bar{r}) \in Z$,

→ THE LOSS FUNCTION CAN TELL HOW GOOD IS h ON (\bar{x}, \bar{r}) :

• $l: \mathcal{H} \times Z \mapsto \mathbb{R}^+$

EXAMPLES:

• 0-1 LOSS: $l(h, (\bar{x}, \bar{r})) = \begin{cases} 1, & \text{IF } h(\bar{x}) \neq \bar{r} \\ 0, & \text{IF } h(\bar{x}) = \bar{r} \end{cases}$

• SQUARED LOSS: $l(h, (\bar{x}, \bar{r})) = (h(\bar{x}) - \bar{r})^2$

• ABSOLUTE VALUE LOSS: $l(h, (\bar{x}, \bar{r})) = |h(\bar{x}) - \bar{r}|$

• COST-SENSITIVE LOSS: $l(h, (\bar{x}, \bar{r})) = C_{h(\bar{x}), \bar{r}}$

/ C IS A $|Y| \times |Y|$ MATRIX

GENERAL PAC LEARNING PROBLEM:

WE WISH TO PA SOLVE: $\min_{h \in H} L_D(h) / L_D(h) := E_{Z \sim D} [l(h, z)]$

$$\frac{1}{m} \sum_i^m l(h, z)$$

- LEARNER KNOWS H, Z AND l , AND RECEIVES ϵ, δ AS PARAMETERS

ACCURACY ↗ CONFIDENCE ↘

- LEARNER CAN DECIDE ON TRAINING SET SIZE m , BASED ON ϵ, δ

- LEARNER DOESN'T KNOW D , BUT CAN SAMPLE $S \sim D^m$

→ USING S , LEARNER OUTPUTS AN HYPOTHESIS $A(S)$

→ WE WANT TO HOLD: $L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon$, WITH $p \geq 1 - \delta$

- A MORE FORMAL DEFINITION:

AN HYPOTHESIS CLASS H IS AGNOSTIC PAC LEARNABLE TO A SET $Z = X \times Y$

AND A LOSS FUNCTION $l : H \times Z \rightarrow \mathbb{R}^+$,

IF ∃ A FUNCTION $m_H : (0, 1)^2 \rightarrow \mathbb{N}$ AND A LEARNING ALGORITHM A WITH THE PROPERTY:

THEORETICAL TIME DISTRIBUTION
↑ OF Z

$\forall \epsilon, \delta \in (0, 1)$, $m \geq m_H(\epsilon, \delta)$, AND A DISTRIBUTION D OVER Z .

$$\left[\begin{array}{c} \text{DATASET } S \in Z^m : L_D(A(S)) \leq \min_{h \in H} L_D(h) + \epsilon \\ \text{ALL POSSIBLE } m\text{-TUPLES IN } Z \\ \text{LEARNER } A \text{ ON } S \\ \text{LOSS ON BEST } h \\ \text{EXPECTED LOSS IS AT MOST } \epsilon\text{-LARGER} \\ \text{THAN THE EXPECTED LOSS OF THE BEST } h \end{array} \right] \geq 1 - \delta$$

CARDESIAN PRODUCT $\hookrightarrow D \times \dots \times D$
EACH $Z \sim D$

• LEARNING VIA UNIFORM CONVERGENCE:

• REPRESENTATIVE SAMPLE:

• A TRAINING SET S IS ϵ -REPRESENTATIVE SAMPLE IF

$$\forall h \in H, |L_s(h) - L_0(h)| \leq \epsilon$$

• LEMMA:

ASSUME THAT A TRAINING SET S IS $\frac{\epsilon}{2}$ -REPRESENTATIVE

\rightarrow ANY OUTPUT OF $\text{ERM}_H(S)$: $h_s \in \arg \min_{h \in H} L_s(h)$

SATISFIES: $L_0(h_s) \leq \min_{h \in H} L_0(h) + \epsilon$

• \rightarrow IF S IS $\frac{\epsilon}{2}$ -REPRESENTATIVE $\rightarrow L_0(h_s) \leq \min_{h \in H} L_0(h) + \epsilon$

• PROOF ($\forall h \in H$): CENTER

$$L_0(h_s) \leq L_s(h_s) + \frac{\epsilon}{2} \stackrel{\text{center}}{\leq} L_s(h) + \frac{\epsilon}{2} \stackrel{L_s(h)}{\leq} L_0(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} \leq L_0(h) + \epsilon$$

• UNIFORM CONVERGENCE: \rightarrow IF VC $\rightarrow H$ CAN LEARN WELL FROM S

H HAS THE UNIFORM CONVERGENCE PROPERTY

IF \exists FUNCTION $m_H^{\text{UC}}: (0, 1)^2 \mapsto \mathbb{N}$ / $\forall \epsilon, S \in (0, 1)$ AND $\forall D$.

$$\mathbb{E}_D \left(\left\{ S \in \mathcal{Z}^m : S \text{ IS } \epsilon\text{-REPRESENTATIVE} \right\} \right) \geq 1 - S$$

• $\rightarrow \exists m_H^{\text{UC}} / \mathbb{E}_D \left(\left\{ S \in \mathcal{Z}^m : S \text{ IS } \epsilon\text{-REPRESENTATIVE} \right\} \right) \geq 1 - S \rightarrow \text{VC}$

• COROLLARY: \rightarrow IF VC $\rightarrow \exists$ BOUND (UPPER) FOR $m_H(\epsilon, S)$

• IF H HAS VC PROPERTY WITH A FUNCTION m_H^{UC}

$\rightarrow H$ IS AGMOTICALLY PAC LEARNABLE WITH $m_H(\epsilon, S) \leq m_H^{\text{UC}}\left(\frac{\epsilon}{2}, S\right)$

• IN THAT CASE, ERM_H ALGORITHM IS A SUCCESSFUL AGMOTIC PAC LEARNER FOR H

- THEOREM (FINITE CLASS \rightarrow AGNOSTIC PAC LEARNABLE) : (NO PROOF, FOR NOW)

LET H BE FINITE AND $l \in [0, 1]$

$\rightarrow H$ IS AGNOSTICALLY PAC LEARNABLE WITH:

$$m_H(\varepsilon, \delta) = \left\lceil \frac{2 \log(2^{|H|}/\delta)}{\varepsilon^2} \right\rceil \rightsquigarrow \left\lceil \dots \right\rceil; \text{ CEILING FUNCTION}$$

$\hookrightarrow \lceil 2.4 \rceil = 3$

- DISCRETIZATION TRICK:

SUPPOSE H PARAMETERIZED BY d NUMBERS / EACH NUMBER REPRESENTED BY b BITS

$$\rightarrow |H| \leq 2^{d \cdot b}$$

$$\rightarrow m_H(\varepsilon, \delta) \leq \left\lceil \frac{2db + 2 \log(2/\delta)}{\varepsilon^2} \right\rceil$$

• LINEAR REGRESSION AND LEAST SQUARE:

- $x \in \mathbb{R}^d, y \in \mathbb{R}$
- $H = \{ \bar{x} \mapsto \langle \bar{w}, \bar{x} \rangle : w \in \mathbb{R}^d \}$

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

- ERM PROBLEM: $\min_{\bar{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m (\langle \bar{w}, \bar{x}_i \rangle - y_i)^2 \rightarrow$ MINIMIZE MSE
OR, EQUIVALENTLY: $\min_{\bar{w} \in \mathbb{R}^d} \left\| \begin{matrix} X^T \\ m \times d \end{matrix} \cdot \begin{matrix} \bar{w} \\ d \times 1 \end{matrix} - \begin{matrix} \bar{y} \\ m \times 1 \end{matrix} \right\|^2$

- GRADIENT AND OPTIMIZATION:

$$\nabla f(\bar{x}) = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right] \quad / f: \mathbb{R}^d \mapsto \mathbb{R}^d$$

\rightarrow IF \bar{x} MINIMIZES $f(\bar{x}) \rightarrow \nabla f(\bar{x}) = (0, \dots, 0)$

- JACOBIAN:

$$\tilde{f}: \mathbb{R}^m \mapsto \mathbb{R}^m$$

$$J_x(\tilde{f}) = \begin{pmatrix} \nabla f_1(\bar{x}) \\ \vdots \\ \nabla f_m \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}, \dots, \frac{\partial f_1}{\partial x_m} \\ \vdots \\ \frac{\partial f_m}{\partial x_1}, \dots, \frac{\partial f_m}{\partial x_m} \end{pmatrix}$$

$$\text{as. } f(x, v) = (\sin x + v, xv) \rightarrow J_x(\tilde{f}) = \begin{pmatrix} \frac{\partial f_1}{\partial x}, \frac{\partial f_1}{\partial v} \\ \frac{\partial f_2}{\partial x}, \frac{\partial f_2}{\partial v} \end{pmatrix} = \begin{pmatrix} \cos x & 1 \\ v & x \end{pmatrix}$$

$$\text{IF } \tilde{f}(\bar{w}) = A \bar{w} \rightarrow J_x(\tilde{f}) = A$$

- CHAIN RULE:

DEFINITION $f: \mathbb{R}^m \mapsto \mathbb{R}^m, g: \mathbb{R}^k \mapsto \mathbb{R}^n$:

$$J_x(f \circ g) = J_{g(x)}(f) \cdot J_x(g)$$

$$\frac{df}{dg} \cdot \frac{dg}{dx}$$

\rightarrow ANALOG OF
 $(f(g(x)))' = f'(g(x)) \cdot g'(x)$

• LEAST SQUARES : \rightarrow PROBLEM: $\min_{\bar{w} \in \mathbb{R}^d} \frac{1}{2} \|X^T \bar{w} - \bar{v}\|^2$

\hookrightarrow SOL: $\bar{w} = (X^T X)^{-1} X^T \bar{v}$ LSP exists

RECALL: WE'D LIKE TO SOLVE THE EPI PROBLEM $\min_{\bar{w} \in \mathbb{R}^d} \frac{1}{2} \|X^T \bar{w} - \bar{v}\|^2$

• LET $\tilde{y}(\bar{w}) = X^T \bar{w} - \bar{v}$, $\tilde{f}(\bar{v}) = \frac{1}{2} \|\bar{v}\|^2 = \frac{1}{2} \sum_{i=1}^m v_i^2$

\hookrightarrow EPI PROBLEM $\hookrightarrow \min_{\bar{w}} \tilde{f}(\tilde{y}(\bar{w}))$ $\hookrightarrow \begin{cases} J_{\bar{w}}(\bar{v}) = X^T \\ J_v(f) = (v_1, \dots, v_m) \end{cases}$

$\rightarrow J_{\bar{w}}(\tilde{f} \circ \bar{v}) = J_{\tilde{y}(\bar{w})}(\tilde{f}) \cdot J_{\bar{w}}(\bar{v}) = \tilde{y}(\bar{w})^T X^T = (X^T \bar{w} - \bar{v})^T X^T$ CHAIN RULE

• OPTIMIZATION \rightarrow SET $J_{\bar{w}}(\tilde{f} \circ \bar{v}) = \bar{0}$

$$\begin{aligned} & (X^T \bar{w} - \bar{v})^T \cdot X^T = 0 \\ & ((X^T \bar{w})^T - \bar{v}^T) X^T = 0 \quad \rightarrow (X^T \bar{w} - \bar{v})^T \cdot X^T = \bar{0}^T \quad \rightarrow X \cdot X^T \frac{\bar{w}}{\bar{v}} = X \bar{v} \\ & (\bar{w}^T X - \bar{v}^T) X^T = 0 \\ & \bar{w}^T X X^T - \bar{v}^T X^T = 0 \quad \rightarrow \text{IF } \exists (X \cdot X^T)^{-1} \quad \rightarrow \bar{w} = (X \cdot X^T)^{-1} \cdot X \bar{v} \\ & \text{UNKNOWN} \end{aligned}$$

• INTERPRETATION AS PROJECTION;

IF \bar{w} IS THE LEAST SQUARE SOLUTION:

$\rightarrow \hat{v} = X^T \bar{w}$ IS THE VECTOR IN $C = \{X^T \bar{w} : \bar{w} \in \mathbb{R}^d\}$ CLOSEST TO \bar{v}

\hookrightarrow PROJECTION OF \bar{v} ONTO C

• WE CAN FIND $\hat{v} = V \cdot V^T \cdot \bar{v}$ $V_{m \times d}$: COLUMNS ARE ORTHONORMAL BASIS OF THE RANGE OF X^T

POLYNOMIAL FITTING :

- 1-DIMENSIONAL POLYNOMIAL FUNCTION OF DEGREE n :

$$p(x) = Q_0 + Q_1 \cdot x + Q_2 \cdot x^2 + \dots + Q_n \cdot x^n$$

- GOAL: given data $S = ((x_1, v_1), \dots, (x_m, v_m)) \rightarrow$ find EBM with respect to the class of polynomials of degree n

- REDUCTION TO LINEAR REGRESSION:

• DEFINE $\psi : \mathbb{R} \mapsto \mathbb{R}^{n+1}$, $\psi(x) = (1, x, x^2, \dots, x^n)$

• DEFINE $\bar{a} = (Q_0, Q_1, \dots, Q_n)$ AND OBSERVE:

$$p(x) = \sum_{i=0}^n Q_i \cdot x^i = \langle \bar{a}, \psi(x) \rangle \rightsquigarrow \text{HALF SPACE / LIN. REGRESSION}$$

\rightarrow TO FIND \bar{a} WE CAN SOLVE LS w.r.t. $\left((\psi(x_1), v_1), \dots, (\psi(x_m), v_m) \right)$

\rightsquigarrow LEAST SQUARES

BIAS - COMPLEXITY TRADE OFF :

• ERROR DECOMPOSITION :

LET $h_s = \text{ERM}_H(S)$. WE CAN DECOMPOSE THE RISK OF h_s AS :

$$L_D(h_s) = \underbrace{\varepsilon_{\text{APP}}}_{\substack{\text{APPROXIMATION} \\ \text{ERROR}}} + \underbrace{\varepsilon_{\text{EST}}}_{\substack{\text{ESTIMATION} \\ \text{ERROR}}}$$

$\hookrightarrow H$

• APPROXIMATION ERROR $\rightarrow \varepsilon_{\text{APP}} = \min_{h \in H} L_D(h)$

2).
CHANGING
POLYNOMIAL DEGREE

\rightarrow IT IS THE RISK WE MAKE DUE TO RESTRICTING TO H \rightarrow 2). FIXING n AS POLYNOMIAL DEGREE

$\uparrow n \rightarrow \downarrow \varepsilon_{\text{APP}}$ \rightarrow IT DOES NOT DEPEND ON S

$\rightarrow \uparrow (\text{complexity of } H / \text{VC dim}) \rightarrow \downarrow \varepsilon_{\text{APP}}$

• ESTIMATION ERROR $\rightarrow \varepsilon_{\text{EST}} = L_D(h_s) - \varepsilon_{\text{APP}}$

\rightarrow IT TURNS HOW GOOD AM I IN EXPANDING H ON MY PROBLEM

\rightarrow IT IS THE RESULT OF L_s BEING ONLY AN ESTIMATE OF L_D

$\rightarrow \uparrow \text{size of } S \rightarrow \downarrow \varepsilon_{\text{EST}}$

$\rightarrow \uparrow \text{complexity of } H \rightarrow \uparrow \varepsilon_{\text{EST}} \because \uparrow \text{complexity} \rightarrow \downarrow \varepsilon_{\text{APP}} \rightarrow \uparrow \varepsilon_{\text{EST}}$

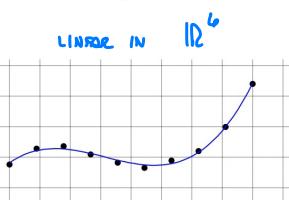
IT CANNOT FIT WELL
OTHER DATA NOT IN TRAINING SET

• HOW TO CHOOSE H ? \rightarrow BIAS - COMPLEXITY TRADE OFF

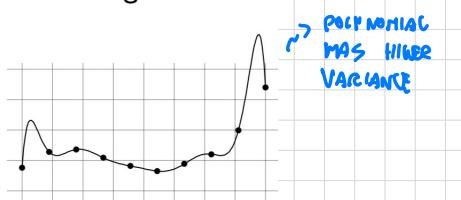
degree 2



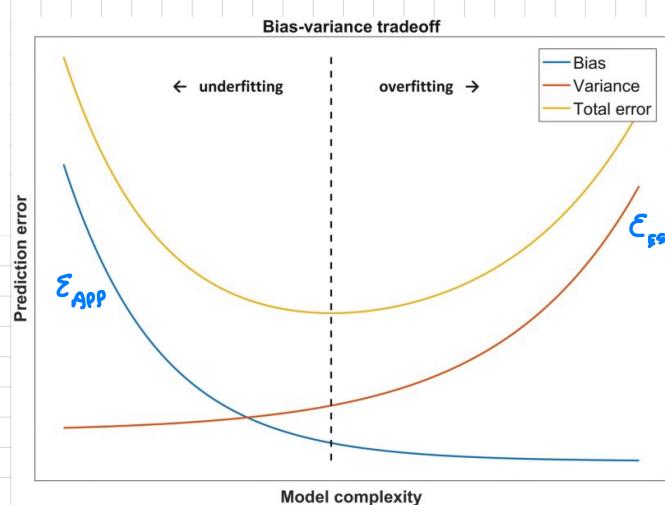
degree 3



degree 10



\hookrightarrow POLYNOMIAL HAS HIGHER VARIANCE



• VALIDATION AND MODEL SELECTION:

• VALIDATION:

SUPPOSE WE ACTUATE USING SOME HYPOTHESIS h , AND WE WANT TO ESTIMATE HOW GOOD IT IS

→ TAKE NEW i.i.d SAMPLE $V = (X_1, Y_1), \dots, (X_{m_V}, Y_{m_V})$ / m_V : # SAMPLE IN VALIDATION

→ IT OUTPUTS $L_V(h)$ AS AN ESTIMATOR OF $L_D(h)$

• USING HOEFDING'S INEQUALITY, IF $\delta \in [0, 1]$:

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

• VALIDATION FOR MODEL SELECTION:

SO,

LET: ○ VALIDATION POINTS, ● TRAINING POINTS

→ DEGREE 3 POLYNOMIAL HAS

MINIMAL VALIDATION ERROR

LET $H = \{h_1, \dots, h_r\}$ BE THE OUTPUT PREDICTION OF APPLYING ERM W.R.T. THE DIFFERENT CLASSES ON S

• LET V BE A NEW VALIDATION SET, CHOOSE $h^* \in \text{ERM}_H(V)$

$$\rightarrow L_D(h^*) \leq \min_{h \in H} L_D(h) + \sqrt{\frac{2\log(2|m_H|/\delta)}{|V|}}$$

• TRAIN - VALIDATION - TEST SETS:

• TRAINING SET: APPLY LEARNING ALGORITHM WITH DIFFERENT

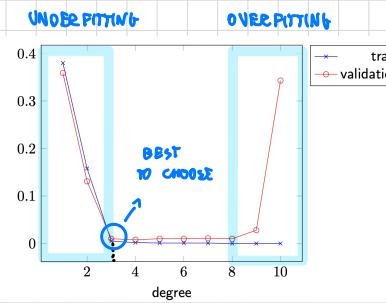
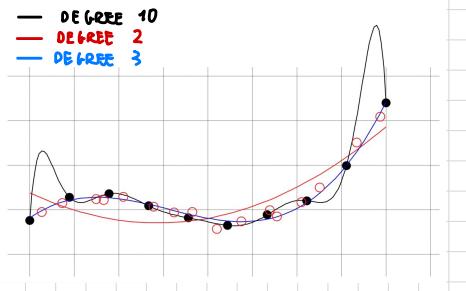
PARAMETERS ON TRAINING SET TO PRODUCE

$$H = \{h_1, \dots, h_n\}$$

• VALIDATION SET: CHOOSE h^* FROM H BASED ON VALIDATION SET \rightsquigarrow TUNING AND SELECTION

• TEST SET: ESTIMATE ERROR OF h^* USING TEST SET

→ K-FOLD IS A GOOD PROCEDURE



SVM AND KERNEL METHODS:

SVM:

HALF-SPACE CLASSIFIER:

$$x \mapsto \text{sign}(\langle \bar{w}, \bar{x} \rangle + b) = \pm 1, \quad \bar{w} \in \mathbb{R}^d$$

MARGIN:

GIVEN THE HYPERPLANE $L = \{v : \langle w, v \rangle + b = 0\}$

SET OF POINTS
OF HYPERPLANE

GIVEN x , DISTANCE OF x TO L :

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

$$\hookrightarrow \text{if } \|w\| = 1 \rightarrow d(x, L) = |\langle \bar{w}, \bar{x} \rangle + b|$$

es. IN \mathbb{R}^3

$$\bar{w} = (a, b, c)$$

PLANE $\tilde{\pi} : ax + by + cz + d = 0$, POINT $P = (x_p, y_p, z_p)$

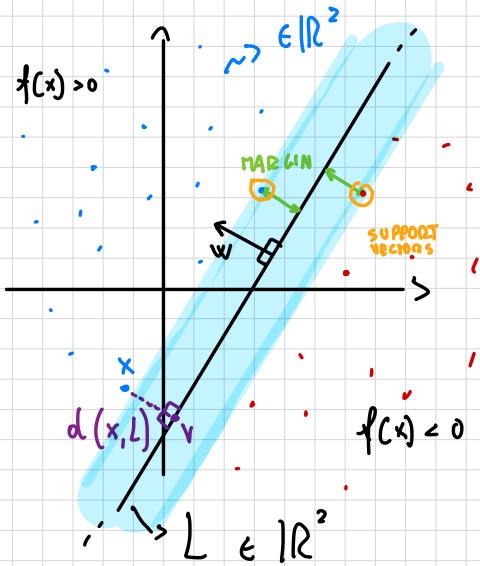
$$\rightarrow d(P, \tilde{\pi}) = \frac{|\langle \bar{w}, \bar{P} \rangle + d|}{\|w\|} = \frac{|ax_p + by_p + cz_p + d|}{\sqrt{a^2 + b^2 + c^2}}$$

MARGIN = $\min_i |\langle \bar{w}, \bar{x}_i \rangle + b| \rightsquigarrow$ CLOSEST DISTANCE FROM L TO x_i, y_i

SUPPORT VECTORS:

A SEPARATING HYPERPLANE IS DEFINED BY $\bar{w}, b / \forall i, y_i \cdot (\langle \bar{w}, \bar{x}_i \rangle + b) > 0$

\rightarrow THE CLOSEST EXAMPLES x_i ARE CALLED SUPPORT VECTORS



$$f(x) : \begin{cases} > 0 \rightarrow \text{CLASS } +1 \\ < 0 \rightarrow \text{CLASS } -1 \end{cases}$$

HARD SVM :

SEEK FOR THE SEPARATING HYPERPLANE WITH LARGEST MARGIN:

$$\text{FIND } \bar{w}, b \text{ TO} \quad \begin{array}{l} \text{MAXIMIZE min distance} \\ \text{argmax}_{(\bar{w}, b) : \|\bar{w}\| = 1} \left(\min_{i \in [m]} \underbrace{\left| \langle \bar{w}, \bar{x}_i \rangle + b \right|}_{\geq 0} \right) \end{array} \quad \begin{array}{l} \text{2 CLASSES SEPARATION} \\ \forall i, \gamma_i (\langle \bar{w}, \bar{x}_i \rangle + b) \geq 0 \end{array}$$

$$\cdot \text{ OR EQUIVALENTLY: } \text{argmax}_{(\bar{w}, b) : \|\bar{w}\| = 1} \left(\min_{i \in [n]} \gamma_i (\langle \bar{w}, \bar{x}_i \rangle + b) \right)$$

• OR EQUIVALENTLY: BEST PLANE: $\langle \bar{w}_0, \bar{x} \rangle + b_0$,

$$(\bar{w}_0, b_0) = \underset{\bar{w}, b}{\text{argmin}} \left(\|\bar{w}\|^2 \right) \quad \forall i, \gamma_i (\langle \bar{w}_0, \bar{x}_i \rangle + b_0) \geq 1$$

$$\rightarrow \text{MARGIN OF } \left(\frac{\bar{w}_0}{\|\bar{w}_0\|}, \frac{b_0}{\|\bar{w}_0\|} \right) = \frac{1}{\|\bar{w}_0\|} \rightarrow \text{MAXIMAL MARGIN}$$

ex.

$$\tilde{w} : \frac{1}{\sqrt{13}} (2x + 3y - 3) = 0 \rightarrow \left(\frac{(2, 3)}{\sqrt{13}}, \frac{-3}{\sqrt{13}} \right)$$

$$\rightarrow \text{MAXIMAL MARGIN} = \frac{1}{\|\bar{w}_0\|} = \frac{1}{\sqrt{13}} \left(= \frac{1}{\sqrt{2^2 + 3^2}} \right)$$

• MARGIN IS SCALE SENSITIVE:

If (\bar{w}, b) separates x_i with margin $\gamma \rightarrow \alpha \cdot x_i$ separate by $\alpha \cdot \gamma$

• MARGIN OF DISTRIBUTION:

D is separable with A (γ, δ) -margin, if $\exists (\bar{w}^*, b^*) / \|\bar{w}^*\| = 1$

AND:

$$D \left(\left\{ (\bar{x}, \gamma) : \|\bar{x}\| \leq \gamma \right\} \cap \left\{ (\bar{x}, \gamma) : \langle \bar{w}^*, \bar{x} \rangle + b^* \geq 1 \right\} \right) = 1$$

\rightarrow THEOREM:

$$\text{if } D \text{ separable with } (\gamma, \delta)-\text{margin} \rightarrow m(\epsilon, \delta) \leq \frac{8}{\epsilon^2} \left(2 \left(\frac{\gamma}{\delta} \right)^2 + \log \left(\frac{2}{\delta} \right) \right)$$

\hookrightarrow UNLIKE VC BOUNDS, $m(\cdot)$ $\frac{8}{\delta}$, NOT ON α

SAMPLE
COMPLEXITY
↑

SOFT-SVM:

- HARD-SVM ASSUMES DATA ARE LINEARLY SEPARABLE

-> IF NOT SEPARABLE -> SOFT-SVM : REGEX CONSTRAINT

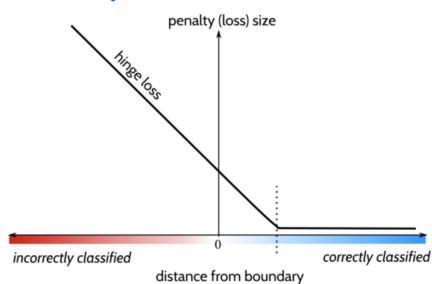
$$\underset{\bar{w}, b, \xi}{\operatorname{argmin}} \left(\lambda \cdot \|\bar{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \quad \text{subject to } \forall i, r_i (\langle \bar{w}, x_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0$$

HINGE

• GIVEN THE HINGE LOSS $\ell^{\text{Hinge}}((\bar{w}, b), (\bar{x}, r)) = \max \{0, 1 - \langle \bar{w}, \bar{x} \rangle - b\}$, IT CAN BE WRITTEN AS REGULARIZED LOSS MINIMIZATION:

$$\underset{\bar{w}, b}{\operatorname{argmin}} \left(\lambda \cdot \|\bar{w}\|^2 + L_s^{\text{Hinge}}((\bar{w}, b)) \right)$$

FOR WHOLE DATASET



• REGULARIZED LOSS MINIMIZATION (RLM):

GIVEN A REGULARIZATION FUNCTION $R: \mathbb{R}^d \mapsto \mathbb{R}$:

$$\text{RLM RULE: } A(S) = \underset{\bar{w}}{\operatorname{argmin}} (L_s(\bar{w}) + R(\bar{w}))$$

• TIKHONOV / RIDGE REGULARIZATION: $R(\bar{w}) = \lambda \|\bar{w}\|^2$

b TERM CAN BE REMOVED IF A NEW FEATURE $x_{d+1} = 1$ IS ADDED

STABILITY:

INFORMALLY, AN ALGORITHM A IS STABLE IF A SMALL CHANGE ON THE INPUTS S

WILL LEAD TO A SMALL CHANGE OF ITS OUTPUTS HYPOTHESIS

• REPLACE 1 SAMPLE: LET $S^{(i)} = (z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m) \quad S = (z_1, \dots, z_n)$

REPLACED SAMPLE POSITION i

ORIGINAL DATA

• ON-AVERAGE-REPLACE-ONE-STABLE DEFINITION:

LET $\varepsilon: \mathbb{N} \mapsto \mathbb{R}$ BE A MONOTONICALLY DECREASING FUNCTION.

• ALGORITHM A IS ON-AVERAGE-REPLACE-ONE-STABLE WITH RATE $\varepsilon(m)$

IF A DISTRIBUTION D :

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

• THEOREM :

IF A IS ON-AVERAGE-REPLACE-ONE-STABLE WITH RATE $\mathcal{E}(m)$:

LOSS BY TRUE DISTRIBUTION D OF DATA

LOSS BY TRAINING DATA

$$\rightarrow \mathbb{E}_{S \sim D^m} [L_0(A(S)) - L_s(A(S))] \leq \mathcal{E}(m)$$

• THEOREM :

TIKHONOV REGULARIZATION AS STABILIZED : $\nabla f(x_1) - \nabla f(x_2) \leq g \cdot |x_1 - x_2|$

LIKE:

ASSUME THAT THE LOSS FUNCTION IS CONVEX AND g -LIPSCHITZ.

\rightarrow RLM RULE, WITH REGULARIZER $\lambda \|w\|^2$, IS ON-AVERAGE-REPLACE-ONE-STABLE

WITH RATE $\mathcal{E}(m) = \frac{2s^2}{\lambda m}$, AND :

RATE $\mathcal{E}(m)$

$$\mathbb{E}_{S \sim D^m} [L_0(A(S)) - L_s(A(S))] \leq \frac{2s^2}{\lambda m}$$

∇f RAPID CHANGES

∇f SMOOTH

• SIMILARLY FOR CONVEX, β -SMOOTH, NON-NEGATIVE LOSS

$$\text{RATE } \mathcal{E}(m) = \frac{48\beta C}{\lambda m} \quad / \quad C : \text{UPPER BOUND ON } \max_z l(\bar{w}, z)$$

• FITTING-STABILITY TRADEOFF:

$$\mathbb{E}_s [L_0(A(S))] = \underbrace{\mathbb{E}_s [L_s(A(S))]}_{\substack{\text{1ST TERM: HOW GOOD } A \text{ FITS } S \\ \text{TRUE RISK}}} + \underbrace{\mathbb{E}_s [L_0(A(S)) - L_s(A(S))]}_{\substack{\text{2ND TERM: OVERFITTING} \\ \text{TRAINING ERROR}}}$$

$\rightarrow \lambda$ CONTROLS TRADEOFF BETWEEN THE 2 TERMS; $\leq \mathcal{E}(m) = \frac{2s^2}{\lambda m}$

$\uparrow \lambda \rightarrow \downarrow \mathcal{E}(m) \rightarrow \downarrow \text{OVERFITTING}$

• FIX SOME ARBITRARY VECTOR \bar{w}^* :

$$L_s(A(S)) \leq L_s(A(S)) + \lambda \|A(S)\|^2 \leq L_s(\bar{w}^*) + \lambda \|\bar{w}^*\|^2$$

TAKING $E[\cdot]$ ON BOTH SIDES $/ E_s[L_s(\bar{w}^*)] = L_0(\bar{w}^*)$

$$\rightarrow \text{WE CAN OBTAIN: } \mathbb{E}_s [L_s(A(S))] \leq L_0(\bar{w}^*) + \lambda \|\bar{w}^*\|^2$$

$$\rightarrow \mathbb{E}_s [L_0(A(S))] \leq L_0(\bar{w}^*) + \lambda \|\bar{w}^*\|^2 + \frac{2s^2}{\lambda m}$$

THE REGULARIZATION PATH:

• RLM rule:

$$\bar{w}(\lambda) = \underset{\bar{w}}{\operatorname{arg\,min}} \left(L_s(\bar{w}) + \lambda \|\bar{w}\|^2 \right)$$

→ WE WANT TO MINIMIZE BOTH $L_s(\bar{w})$ AND $\|\bar{w}\|^2$:

• $\bar{w}(\lambda=0)$: MINIMIZE JUST $L_s(\bar{w})$

• $\bar{w}(\lambda=\infty)$: MINIMIZING $\|\bar{w}\|^2$ AS AN HIGHER IMPACT: $L_s(\bar{w}) \xrightarrow{\lambda \|\bar{w}\| \gg 0} 0$

• HOW TO CHOOSE λ ?

→ BOUND MINIMIZATION: choose λ according to the bounds on $L_s(\bar{w})$ USUALLY FAR FROM OPTIMAL, AS THE BOUND IS WORST CASE

→ VALIDATION: CALCULATE SEVERAL PARETO OPTIMAL POINTS ON THE REGULARIZATION PATH (BY VARYING λ) AND USE VALIDATION SET TO CHOOSE THE BEST ONE

SAMPLE COMPLEXITY OF SOFT-SVM.

• SOFT-SVM \equiv RLM

$$g = \|\bar{x}\|$$

↑

• HINGE-LOSS: $\bar{w} \mapsto \max \{0, 1 - \gamma \langle \bar{w}, \bar{x} \rangle\}$ IS $\|\bar{x}\|$ -LIPSHITZ:

$$\|\bar{x}\| \text{-LIPSHITZ} \rightarrow |\max \{0, 1 - \gamma \langle \bar{w}, \bar{x} \rangle\} - \max \{0, 1 - \gamma \langle \bar{w}', \bar{x} \rangle\}| \leq \|\bar{x}\| \cdot \|\bar{w} - \bar{w}'\|$$

• WITH MIN-MAX SCALER: $\frac{\bar{x} - \min(\bar{x})}{\max(\bar{x}) - \min(\bar{x})} \rightarrow \text{DATA} \subset [0,1]^d \rightarrow \exists g / \|\bar{x}\| < g$

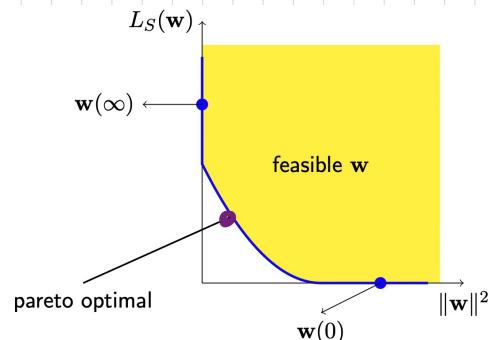
• ASSUME 0 SUCH THAT $\|\bar{x}\| < g$ WITH PROGRAMM = 1

→ WE OBTAIN A CONVEX-LIPSHITZ LOSS, AND $\nabla \bar{u}$:

$$\underset{s \sim 0^m}{\mathbb{E}} [L_0^{\text{HINGE}}(A(s))] \leq L_0^{\text{HINGE}}(\bar{u}) + \lambda \|\bar{u}\|^2 + \frac{2s^2}{\lambda m}$$

• IF WE SET $\lambda = \frac{1}{B} \sqrt{\frac{2g^2}{cm}}$ CONST., AND SINCE L_0^{HINGE} IS AN UPPER BOUND FOR 0-1 LOSS:

$$\rightarrow \underset{s \sim 0^m}{\mathbb{E}} [L_0^{0-1}(A(s))] \leq \min_{\bar{w}: \|\bar{w}\| \leq B} \{L_0^{\text{HINGE}}(\bar{w})\} + \sqrt{\frac{8s^2 B^2}{m}}$$



MARGIN / NORM VS DIMENSIONALITY :

- VC dimension of learning half spaces (\hookrightarrow) dimension of \sim
- $\rightarrow \uparrow \text{dim} \rightarrow \uparrow \text{sample complexity}$
- for SVMs : sample complexity (\hookrightarrow) $\left(\frac{\gamma}{\delta}\right)^2 = (\gamma^2 B^2)$

(SEE GPT BOOK MEHKS
FOR EXAMPLES -> NOTION)

SGD for solving Soft-SVM

(NOT COVERED?)

goal: Solve $\operatorname{argmin}_{\mathbf{w}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\} \right)$

parameter: T

initialize: $\boldsymbol{\theta}^{(1)} = \mathbf{0}$

for $t = 1, \dots, T$

Let $\mathbf{w}^{(t)} = \frac{1}{\lambda t} \boldsymbol{\theta}^{(t)}$

Choose i uniformly at random from $[m]$

If ($y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle < 1$)

Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + y_i \mathbf{x}_i$

Else

Set $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)}$

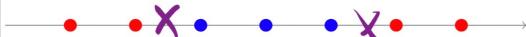
output: $\bar{\mathbf{w}} = \frac{1}{T} \sum_{t=1}^T \mathbf{w}^{(t)}$

KERNELS :

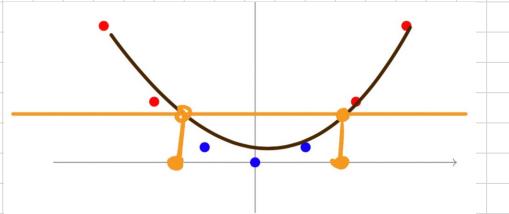
EMBEDDING INTO FEATURE SPACE:

ex.

→ THIS SAMPLE IN \mathbb{R}^n IS NOT SEPARABLE BY HYPERSPLANES



BUT...



GENERAL APPROACH:

1. DEFINE $\psi: X \rightarrow F$ / $F \sim$ SUBSET OF HILBERT SPACE: GENERALIZATION OF EUCLIDEAN SPACE
2. TRAIN A HYPERSPACE OVER $(\psi(\bar{x}_1), v_1), \dots, (\psi(\bar{x}_m), v_m)$

POLYNOMIAL MAPPINGS:

DEGREE K

$$\cdot p(x) = \sum_{j=0}^k w_j x^j = \langle \bar{w}, \psi(x) \rangle \quad \psi(x) = (1, x, x^2, \dots, x^k)$$

more generally, a degree K multivariate polynomial from $\mathbb{R}^n \rightarrow \mathbb{R}$:

$$p(\bar{x}) = \sum_{\substack{j \in [n]^r \\ r \leq k}} w_j \cdot \prod_{i=1}^r x_{j_i} = \langle \bar{w}, \psi(\bar{x}) \rangle \quad \psi: \mathbb{R}^n \rightarrow \mathbb{R}^d$$

ex. $J \in [n]^r$, $n=3$, $r=2$:

$$\rightarrow [3]^2 \rightarrow [3] \times [3] = \begin{matrix} x_1 x_2 \\ x_2 x_1 \end{matrix} (2, 1), (2, 2), (2, 3), \\ (3, 1), (3, 2), (3, 3)$$

$$\rightarrow (x_1, x_2, x_3) \mapsto (x_1^2, x_2^2, x_3^2, x_1 x_2, x_2 x_3, x_1 x_3, x_1, x_2, x_3)$$

KERNEL TRICK:

IT CAN HELP IN avoid to EXPLICITLY COMPUTE $f(x), f(y)$ AND USE THE SHORT CUT $K(\bar{x}_1, \bar{x}_2)$

$$K(\bar{x}_1, \bar{x}_2) = \langle \psi(\bar{x}_1), \psi(\bar{x}_2) \rangle$$

→ SOMETIMES CALCULATING $K(\cdot)$ IS EASIER AND MORE EFFICIENT

ex.

Simple Example: $x = (x_1, x_2, x_3); y = (y_1, y_2, y_3)$. Then for the function $f(x) = (x_1 x_1, x_1 x_2, x_1 x_3, x_2 x_1, x_2 x_2, x_2 x_3, x_3 x_1, x_3 x_2, x_3 x_3)$, the kernel is $K(x, y) = (\langle x, y \rangle)^2$.

EXPLICIT CALCULUS

Let's plug in some numbers to make this more intuitive: suppose $x = (1, 2, 3); y = (4, 5, 6)$. Then:
 $f(x) = (1, 2, 3, 2, 4, 6, 3, 6, 9)$
 $f(y) = (16, 20, 24, 20, 25, 30, 24, 30, 36)$
 $\langle f(x), f(y) \rangle = 16 + 40 + 72 + 40 + 100 + 180 + 72 + 180 + 324 = 1024$

USING $K(x, y)$

Now let us use the kernel instead:
 $K(x, y) = (4 + 10 + 18)^2 = 32^2 = 1024$

Same result, but this calculation is so much easier.

THE REPRESENTER THEOREM:

CONSIDER ANY LEARNING RULE, $f : \mathbb{R}^m \rightarrow \mathbb{R}$, OF THE FORM :

$$\bar{w}^* = \underset{\bar{w}}{\operatorname{arg\min}} \left\{ f \left(\langle \bar{w}, \psi(\bar{x}_1) \rangle, \dots, \langle \bar{w}, \psi(\bar{x}_m) \rangle \right) + \lambda \|\bar{w}\|^2 \right\}$$

$$\rightarrow \exists \bar{\alpha} \in \mathbb{R}^m / \bar{w}^* = \sum_{i=1}^m \alpha_i \cdot \psi(\bar{x}_i) \rightarrow \begin{matrix} \bar{w}^* \text{ IS A LINEAR COMBINATION} \\ \text{OF } \psi(\bar{x}_i) \end{matrix}$$

IMPLICATIONS OF THEOREM:

GRAM MATRIX \rightarrow POSITIVE ($\lambda_i > 0$), SEMI-DEFINED ($\lambda_i \geq 0$)

• DENOTE BY G THE MATRIX / $G_{i,j} = \langle \psi(\bar{x}_i), \psi(\bar{x}_j) \rangle, \forall i$:

REPR. THEORY

$$\rightarrow \langle \bar{w}, \psi(\bar{x}_i) \rangle = \langle \sum_j \alpha_j \psi(\bar{x}_j), \psi(\bar{x}_i) \rangle =$$

$$\langle \bar{w}, \psi(\bar{x}_i) \rangle = (G\bar{\alpha})_i = \sum_{j=1}^m \alpha_j \cdot \langle \psi(\bar{x}_j), \psi(\bar{x}_i) \rangle = \underbrace{(G \cdot \bar{\alpha})_i}_{\in \mathbb{R}}$$

or. $\{x_1, x_2\}$

, $\|\psi(\bar{x}_i)\|^2 \rightarrow$ ANALOG FOR $i, j = 2$

$$G = \begin{pmatrix} \langle \psi(\bar{x}_1), \psi(\bar{x}_1) \rangle & \langle \psi(\bar{x}_1), \psi(\bar{x}_2) \rangle \\ \langle \psi(\bar{x}_2), \psi(\bar{x}_1) \rangle & \langle \psi(\bar{x}_2), \psi(\bar{x}_2) \rangle \end{pmatrix} = \begin{pmatrix} G_{11} & G_{12} \\ G_{21} & G_{22} \end{pmatrix}$$

$= \langle \bar{w}, \bar{w} \rangle$

$$\rightarrow \|\bar{w}\|^2 = \langle \sum_j \alpha_j \psi(\bar{x}_j), \sum_j \alpha_j \psi(\bar{x}_j) \rangle =$$

$$\| \bar{w} \|^2 = \bar{\alpha}^T G \bar{\alpha} = \sum_{i,j=1}^m \alpha_i \cdot \alpha_j \cdot \langle \psi(\bar{x}_i), \psi(\bar{x}_j) \rangle = \bar{\alpha}^T \cdot G \cdot \bar{\alpha}$$

$\left(= \sum_{i=1}^m \alpha_i G_{i,i} + 2 \sum_{i < j} \alpha_i \alpha_j G_{i,j} \right) \rightarrow$ QUADRATIC FORM

\rightarrow SINCE $\langle \bar{w}, \psi(\bar{x}) \rangle = G \cdot \bar{\alpha}$, WE CAN OPTIMIZE OVER $\bar{\alpha}$:

$$\underset{\bar{\alpha} \in \mathbb{R}^m}{\operatorname{arg\min}} \left\{ f(G \cdot \bar{\alpha}) + \lambda \bar{\alpha}^T G \bar{\alpha} \right\}$$

• N.B.: G ONLY DEPENDS ON $\langle \cdot, \cdot \rangle \rightarrow$ IT CAN BE COMPUTED USING $K(\cdot, \cdot)$ ALONE

\rightarrow SUPPOSE WE FOUND $\bar{\alpha}$, THEN GIVEN A NEW INSTANCE \bar{x} :

\bar{x} NEW INSTANCE

$$\langle \bar{w}, \psi(\bar{x}) \rangle = \langle \sum_j \psi(\bar{x}_j), \psi(\bar{x}) \rangle = \sum_j \langle \psi(\bar{x}_j), \psi(\bar{x}) \rangle =$$

$$= \sum_j K(\bar{x}_j, \bar{x}) \rightarrow$$
 TRAINING AND PREDICTIONS CAN BE DONE USING $K(\cdot, \cdot)$ ALONE

• REPRESENTER THEOREM for SVM:

$$\text{soft-SVM} : \min_{\bar{\alpha} \in \mathbb{R}^m} \left\{ \underbrace{\lambda \|\bar{\alpha}\|^2}_{\text{Hinge}} + \frac{1}{m} \sum_{i=1}^m \max \{0, 1 - r_i (\bar{G} \bar{\alpha})_i\} \right\}$$

$$\text{hard-SVM} : \min_{\bar{\alpha} \in \mathbb{R}^m} \left\{ \underbrace{\bar{\alpha}^\top (\bar{G} \bar{\alpha})}_{\text{Margin}} \right\} / \forall i, r_i (\bar{G} \bar{\alpha})_i \geq 1$$

• POLYNOMIAL KERNELS:

$$K \text{ degree polynomial kernel: } K(\bar{x}_1, \bar{x}_2) = (1 + \langle \bar{x}_1, \bar{x}_2 \rangle)^k$$

so, $k=2$

$$K(\bar{x}_1, \bar{x}_2) = (1 + \langle \bar{x}_1, \bar{x}_2 \rangle)^2 = \left(1 + \underbrace{\frac{\bar{x}_{1,1} \cdot \bar{x}_{2,1}}{a}}_a + \underbrace{\bar{x}_{1,2} \cdot \bar{x}_{2,2}}_b \right)^2 = \\ \left(= 1^2 + a^2 + b^2 + 2 \cdot 1 \cdot a + 2 \cdot 1 \cdot b + 2 \cdot a \cdot b = 1 + a^2 + b^2 + 2a \cdot 2b + 2ab \right)$$

$$= 1 + (\bar{x}_{1,1} \cdot \bar{x}_{2,1})^2 + (\bar{x}_{1,2} \cdot \bar{x}_{2,2})^2 + 2 \cdot \bar{x}_{1,1} \bar{x}_{2,1} + 2 \cdot \bar{x}_{1,2} \bar{x}_{2,2} + 2 \cdot \bar{x}_{1,1} \bar{x}_{2,1} \cdot \bar{x}_{1,2} \bar{x}_{2,2}$$

• N.B.: calculating $K(\bar{x}_1, \bar{x}_2)$ takes $O(n)$ time,

while dimension of $\psi(\bar{x})$ is $n^k \rightarrow n$: # samples

• GAUSSIAN KERNELS (RBF):

$$K(\bar{x}_1, \bar{x}_2) = e^{-\frac{\|\bar{x}_1 - \bar{x}_2\|^2}{2\sigma^2}} \rightarrow \text{IT CAN LEARN ANY POLYNOMIAL}$$

• LET ORIGINAL INSTANCE SPACE BE \mathbb{R} , MAPPING ψ / $\exists \psi(x)_n = \frac{1}{\sqrt{n!}} e^{-\frac{x^2}{2}} \cdot x^n$

$$\rightarrow \langle \psi(x_1), \psi(x_2) \rangle = \sum_{n=0}^{\infty} \left(\frac{1}{\sqrt{n!}} e^{-\frac{x_1^2}{2}} \cdot x_1^n \right) \left(\frac{1}{\sqrt{n!}} e^{-\frac{x_2^2}{2}} \cdot x_2^n \right) = \\ = e^{-\frac{x_1^2 + x_2^2}{2}} \cdot \sum_{n=0}^{\infty} \frac{(x_1 x_2)^n}{n!} = e^{-\frac{(x_1 - x_2)^2}{2}}$$

• LEMMA (MERCER'S CONDITION):

A SYMMETRIC FUNCTION $K: X \times X \mapsto \mathbb{R}$ IMPLEMENTS AN INNER PRODUCT

IN SOME HILBERT SPACE $\Leftrightarrow \forall \bar{x}_1, \dots, \bar{x}_m : G_{i,j} = K(\bar{x}_i, \bar{x}_j)$ IS A POSITIVE, SEMI-DEFINITE MATRIX

Fritz John Condition:

THEOREM 15.8 Let \mathbf{w}_0 be as defined in Equation (15.3) and let $I = \{i : |\langle \mathbf{w}_0, \mathbf{x}_i \rangle| = 1\}$. Then, there exist coefficients $\alpha_1, \dots, \alpha_m$ such that

$$\mathbf{w}_0 = \sum_{i \in I} \alpha_i \mathbf{x}_i.$$

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

The proof of this theorem follows by applying the following lemma to Equation (15.3).

LEMMA 15.9 (Fritz John) Suppose that

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w}) \quad \text{s.t. } \forall i \in [m], g_i(\mathbf{w}) \leq 0,$$

where f, g_1, \dots, g_m are differentiable. Then, there exists $\boldsymbol{\alpha} \in \mathbb{R}^m$ such that $\nabla f(\mathbf{w}^*) + \sum_{i \in I} \alpha_i \nabla g_i(\mathbf{w}^*) = \mathbf{0}$, where $I = \{i : g_i(\mathbf{w}^*) = 0\}$.

DUALITY:

• FOR HARD-SVM, FOR EXAMPLE:

$$\min_{\bar{\mathbf{w}}} \|\bar{\mathbf{w}}\|^2 \quad / \quad \forall i, \gamma_i \leq \bar{\mathbf{w}}, \bar{\mathbf{x}}_i \geq 1$$

$$\xrightarrow[\substack{\text{IT CAN BE} \\ \text{WRITTEN AS}}]{} \min_{\bar{\mathbf{w}}} \left\{ \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^m \\ / \bar{\mathbf{z}} \geq 0}} \left\{ \frac{1}{2} \|\bar{\mathbf{w}}\|^2 + \sum_{i=1}^m \alpha_i (\gamma_i - \gamma_i \leq \bar{\mathbf{w}}, \bar{\mathbf{x}}_i) \right\} \right\}$$

LAGRANGIAN
MULTIPLIER

CONDITION

→ WE CAN OBTAIN THE WEAK DUALITY INEQUALITY:

$$\min_{\bar{\mathbf{w}}} \left\{ \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^m \\ / \bar{\mathbf{z}} \geq 0}} \left\{ \frac{1}{2} \|\bar{\mathbf{w}}\|^2 + \sum_{i=1}^m \alpha_i (\gamma_i - \gamma_i \leq \bar{\mathbf{w}}, \bar{\mathbf{x}}_i) \right\} \right\} \geq \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^m \\ / \bar{\mathbf{z}} \geq 0}} \left\{ \min_{\bar{\mathbf{w}}} \left\{ \frac{1}{2} \|\bar{\mathbf{w}}\|^2 + \sum_{i=1}^m \alpha_i (\gamma_i - \gamma_i \leq \bar{\mathbf{w}}, \bar{\mathbf{x}}_i) \right\} \right\}$$

• WE CAN SOLVE ANOTHERLY THE OPTIMIZATION $\min \{ \dots \}$:

$$\xrightarrow{\text{SOLUTION:}} \bar{\mathbf{w}} = \sum_{i=1}^m d_i \gamma_i \bar{\mathbf{x}}_i$$

→ WHICH LEADS TO: → PLUGGING IT BACK

$$\max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^m \\ / \bar{\mathbf{z}} \geq 0}} \left\{ \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i \gamma_i \bar{\mathbf{x}}_i \right\|^2 + \sum_{i=1}^m d_i (\gamma_i - \gamma_i \leq \sum_j d_j \gamma_j \bar{\mathbf{x}}_j, \bar{\mathbf{x}}_i) \right\}$$

CURSE OF DIMENSIONALITY:

GIVEN n IN THE SPACE DIMENSION: $X \in \mathbb{R}^n \rightarrow \text{as. } n=3 : \bar{x} = (x_1, x_2, x_3)$

① $\uparrow n \rightarrow \uparrow \# \text{ OPERATIONS REQUIRED TO COMPUTE DISTANCES:}$

$$d(\bar{x}, \bar{v}) = \|\bar{x} - \bar{v}\| = \sqrt{\sum_i (x_i - v_i)^2} / \bar{x}, \bar{v} \in \mathbb{R}^n$$

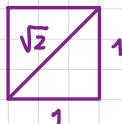
as. $n=3 : d(\bar{x}, \bar{v}) = ((x_1 - v_1)^2 + (x_2 - v_2)^2 + (x_3 - v_3)^2) \rightarrow 3 \text{ SUMS}$

as. $n=100 : d(\bar{x}, \bar{v}) = ((x_1 - v_1)^2 + \dots + (x_{100} - v_{100})^2) \rightarrow 100 \text{ SUMS}$

② $\uparrow n \rightarrow \uparrow \text{ DISTANCES LENGTH:}$

$$Q_d = [0, 1]^d, \text{ MAX DISTANCE} = \| (1, \dots, 1) - (0, \dots, 0) \| = \sqrt{d}$$

as. $d=2$



as. $d=3$



....

as. $d=100$

$$\rightarrow \sqrt{100}$$

→ ① AND ② CAN CAUSE PROBLEMS IN MANY MACHINE LEARNING

PROBLEMS: DISTANCES FOR K-NN, ESTIMATION IN REGRESSIONS, ETC.

$$\text{so. } \mathbb{R}^p : p=2 \rightarrow s^{\frac{1}{2}} = \sqrt{s}$$

HIGH DIMENSIONAL SPACES ARE EMPTY:

ASSUME DATA LIVES IN $[0, 1]^p$.

TO CAPTURE A NEIGHBOURHOOD WHICH IS A FRACTION s OF THE HYPERCUBE VOLUME, EDGE LENGTH = $s^{1/p}$

10% OF VOLUME

$$\cdot s = 0.1, p = 10 \rightarrow s^{1/p} = 0.80$$

$\downarrow /10$

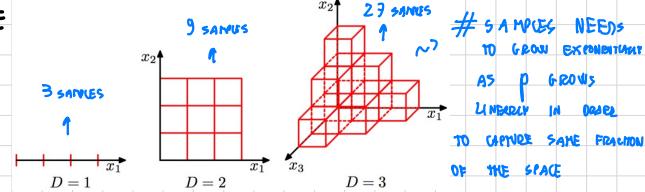
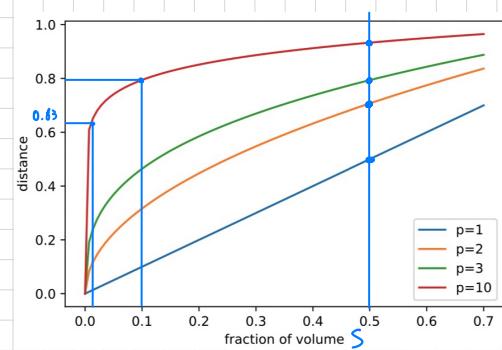
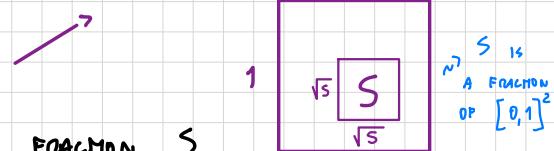
$$\cdot s = 0.01, p = 10 \rightarrow s^{1/p} = 0.63$$

1% OF VOLUME

· IN ORDER TO CAPTURE HIGHER s

→ AS $\uparrow p$, A "LONGER" EDGE

IS NEEDED, WHICH GROWS EXPONENTIALLY



NEAREST NEIGHBOURS :

Given 2 R.V. $X, Y \sim U([0, 1]^p)$

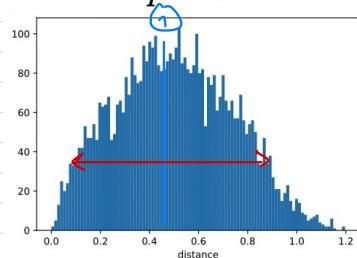
DIMENSIONS

• MEAN SQUARE DISTANCE (MSD) = $\|X - Y\|^2$:

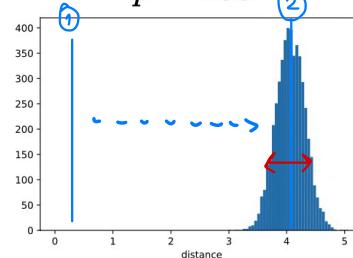
- $E[MSD] = \alpha = \frac{1}{6} p$
- $STD[MSD] = \beta \approx 0.2\sqrt{p}$

$$\left/ \begin{array}{l} \text{SCALED STD:} \\ \frac{\beta}{\alpha} = \frac{\sqrt{p}}{p} \approx p^{-\frac{1}{2}} \end{array} \right.$$

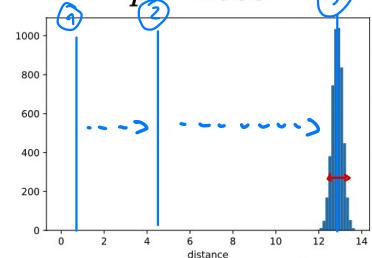
$p = 2$



$p = 100$



$p = 1000$



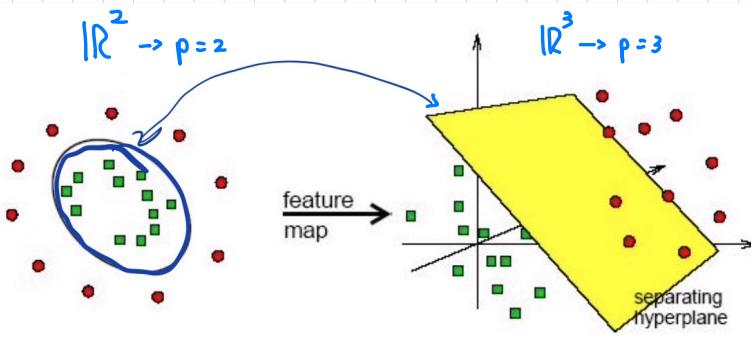
• $\uparrow p \rightarrow \uparrow E[MSD], \downarrow STD[MSD]$

$\rightarrow \uparrow E[MSD]$: ALL POINTS TEND TO BE AT SMALLER DISTANCES FROM EACH OTHER

$\rightarrow \downarrow STD[MSD]$: ALL POINTS TEND TO HAVE ALMOST THE SAME DISTANCE FROM EACH OTHER

\rightarrow THE NOTION OF NN VANISHES : DIFFICULTY TO HAVE MEANINGFUL DISTANCES AND USE THEM

\rightarrow BUT $p \gg 0$ CAN BE USEFUL TO SEPARATE GROUPS BETTER;



complex in low dimensions

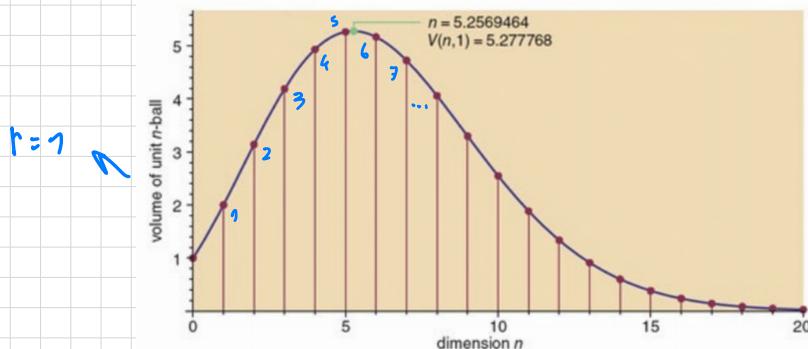
simple in higher dimensions

CONCENTRATION PHENOMENA:

VOLUME OF HYPER-BALL, RADIUS r :

$$V_p(r) = r^p \frac{\pi^{p/2}}{\Gamma(\frac{p}{2} + 1)}$$

GAMMA FUNCTION:
 $\Gamma(n+1) = n!, n \in \mathbb{N}$, $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt, z \in \mathbb{C}$

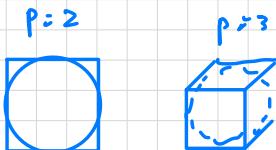


→ IF YOU WANT TO COVER $[0,1]^p$: \rightsquigarrow

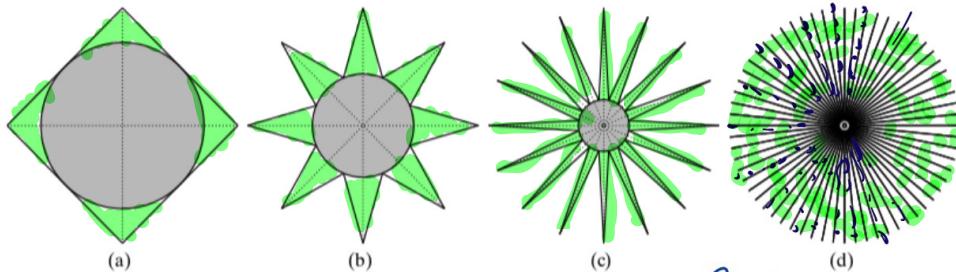
→ STRANGE PHENOMENA:

APTER $n = 5, 25 \dots$

VOLUME STARTS DECREASING!



MOST OF POINTS



ASSUME NOW DRAW n SAMPLES WITH $\sim V$

→ MOST SAMPLES WILL BE IN THE CORNERS OF HYPERCUBE.

PROBABILIT Y FOR $X \sim V([0,1]^p)$ BELONGS TO SMALL

BETWEEN $r = 0.9$ AND $r \geq 1$; \rightsquigarrow AS r_p , IT'S MORE

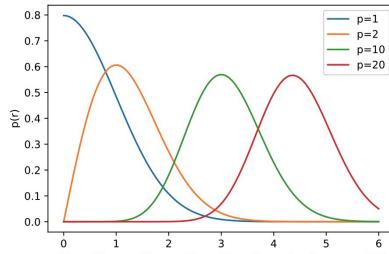
$$P(X \in S_{0.9}(p)) = 1 - 0.9^p \xrightarrow[p \rightarrow \infty]{} 1$$

SURE THAT THE POINT
WILL BE IN $0.9 < r \leq 1$

FOR

GAUSSIAN

DISTRIBUTION:



Most of the mass of a Gaussian distribution is located in areas where the density is extremely small compared to its maximum value.



DIMENSIONALITY REDUCTION:

DEF.: TAKING DATA IN HIGH DIMENSIONAL SPACE AND MAPPING IT INTO A LOW DIMENSIONAL SPACE

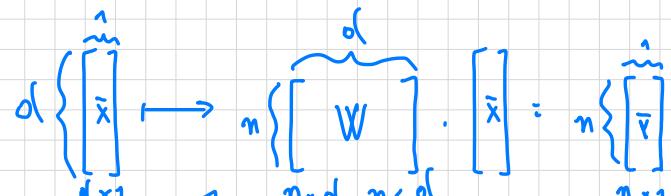
→ REDUCES TRAINING AND TESTING TIME

→ REDUCES ESTIMATION ERROR

→ INTERPRETABILITY OF DATA

LINEAR DIMENSIONALITY REDUCTION:

$$\bar{x} \rightarrow W\bar{x} = \bar{v} / W \in \mathbb{R}^{n,d}, n > d$$



$$O(nd^2 + d^3)$$

PCA (PRINCIPAL COMPONENT ANALYSIS):

$$\bar{x} \xrightarrow{\text{W}} W\bar{x} = \bar{v} \quad / \quad \bar{x} = [x_1, x_2, \dots, x_d]^T$$

• NATURAL QUESTION: WE WANT TO BE ABLE TO APPROXIMATELY RECOVER \bar{x} FROM $\bar{v} = W\bar{x}$

→ IN PCA: $\tilde{x} = \bigcup_{d \times n} \bar{v} = \bigcup_{d \times n \text{ and}} W\bar{x}$ LINEAR RECOVERY

↳ REQUIRES "APPROXIMATE RECOVERY", SOLVE:

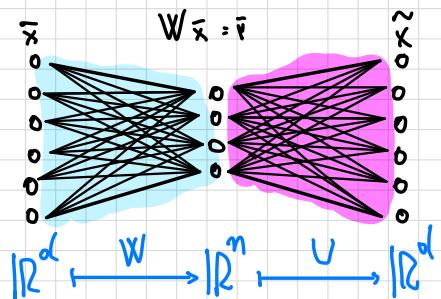
$$\underset{W \in \mathbb{R}^{n,d}, V \in \mathbb{R}^{d,n}}{\operatorname{arg\,min}} \left\{ \sum_{i=1}^m \| \bar{x}_i - \bigcup_{\substack{\text{ORIGINAL FEATURES} \\ \text{RECOVERED FEATURES}}} W\bar{x}_i \|_2^2 \right\} \rightarrow m: \# \text{SAMPLES}$$

• IT IS A LOSSY PROCESS: $\tilde{x} \neq \bar{x}$

$$V \cdot W \approx I_d$$

$$\bar{x} \xrightarrow[\substack{d \times 1 \\ \text{ENCODING}}]{W} W\bar{x} = \bar{v} \xrightarrow[U]{\substack{n \times d \\ \text{DECODING}}} U \cdot \bar{v} = \tilde{x}$$

• IT CAN BE SEEN AS A MLP:



• SOLVING THE PCA PROBLEM:

$$\underset{W \in \mathbb{R}^{n,d}, V \in \mathbb{R}^{d,n}}{\operatorname{arg\,min}} \left\{ \sum_{i=1}^m \| \bar{x}_i - V W \bar{x}_i \|_2^2 \right\}$$

EIGENVALUE λ OF A : $\det(A - \lambda I)$
 EIGENVECTORS \bar{v} OF A : $(A - \lambda I) \cdot \bar{v} = 0$

CORRESPONDING EIGENVALUE OF \bar{v}

• THEOREM: $\underbrace{\text{IF } \mu = 0: \text{ LET } A = \sum_{i=1}^m \bar{x}_i \bar{x}_i^\top}_{\text{A} \in \mathbb{R}^{d,d}}$, $\underbrace{\text{LET } \bar{u}_1, \dots, \bar{u}_m}_{\text{BE THE}} \text{ BE THE LEADING EIGENVECTORS OF } A$:

$$U = \begin{bmatrix} \bar{u}_1 & \bar{u}_2 & \dots & \bar{u}_m \end{bmatrix}$$

→ SOLUTION OF PCA: $\left\{ \begin{array}{l} \text{COLUMNS OF } U \text{ SET TO } \bar{u}_1, \dots, \bar{u}_m \\ W = U^\top \end{array} \right.$

$$(C \cdot D)^\top = D^\top C$$

from (i) \rightarrow (i) PROPERTY: $A^\top = \left(\sum_i \bar{x}_i \bar{x}_i^\top \right)^\top = \sum_i (\bar{x}_i \cdot \bar{x}_i^\top)^\top = \sum_i \bar{x}_i \cdot \bar{x}_i^\top = A \rightarrow A \text{ symmetric}$

(ii) PROPERTY: IF A symmetric $\rightarrow \lambda_i \in \mathbb{R}$

PROPERTY: $\lambda_1 \geq \dots \geq \lambda_d \geq 0 \rightarrow A$ IS POSITIVE SEMI-DEFINITE

• N.B.: $\frac{1}{m-1} (\bar{x} - \mu)^\top (\bar{x} - \mu) = \sum = \text{COVARIANCE MATRIX}$

• N.B.: $\| \bar{x} - V W^\top \bar{x} \|_2^2 = \dots = \| \bar{x} \|_2^2 - \text{trace} \left(W^\top \bar{x} \bar{x}^\top W \right)$

→ EQUIVALENT PCA PROBLEM:

$$\underset{V \in \mathbb{R}^{d,n}, V^\top V = I}{\operatorname{arg\,max}} \left\{ \text{trace} \left(V^\top \left(\sum_{i=1}^m \bar{x}_i \bar{x}_i^\top \right) V \right) \right\}$$

→ SOLUTION IS TO SET V TO BE THE LEADING EIGENVECTORS OF $A = \sum_{i=1}^m \bar{x}_i \bar{x}_i^\top$

• PCA AND SVD :

<https://stats.stackexchange.com/questions/134282/relationship-between-svd-and-pca-how-to-use-svd-to-perform-pca>

COVARIANCE MATRIX



- PCA CAN BE PERFORMED BY FINDING EIGENVECTORS OF SVD ON $X \sim$ DATA MATRIX

RECALL

$$\text{SVD } (X) = U \cdot S \cdot V^T$$

$$XX^T = []$$

$$(X^T X) = []$$

/

U: COLUMNS ARE EIGENVECTORS

OF $X \cdot X^T$

$$\lambda_i = \sqrt{\lambda_i} \text{ or } X^T X$$

S: DIAGONAL MATRIX, $S_{i,i} \geq 0$

V: ROWS ARE EIGENVECTORS OF $X^T X$

$$X \xrightarrow{\mu} X - \bar{\mu}$$

- IN PCA: $\Sigma = \frac{1}{n-1} X^T X$, ASSUMING DATA IN X ARE ZERO-CENTERED

- Σ IS SYMMETRIC $\rightarrow \Sigma$ CAN BE DIAGONALIZED. $\Sigma = V \cdot L \cdot V^T$

WHERE:

EIGENVECTOR \bar{V}_i

EIGENVALUE λ_i

$$V = \begin{bmatrix} | & | & | \\ \bar{V}_1 & \dots & \bar{V}_i & \dots & \bar{V}_d \\ | & | & | \end{bmatrix}$$

$$L = \begin{bmatrix} \lambda_1 & & & 0 \\ & \ddots & \lambda_i & \\ 0 & & \ddots & \lambda_d \end{bmatrix} / \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

DIAGONAL MATRIX

- \bar{V}_i : PRINCIPAL AXES / PROJECTION OF DATA ON $\bar{V}_i \rightarrow PC_i$

$$\hookrightarrow PC_i = [X] \cdot [\bar{V}_i]$$

$$\hookrightarrow X_{i,PCA} = [x_i] \cdot [V]$$

COORDINATES OF i-TH POINTS
ON NEW PC SCALE

- RELATION BETWEEN SVD AND PCA:

$$\Sigma = \frac{1}{n-1} X^T X = \frac{1}{n-1} VSU^T \cdot USV^T = V \cdot \frac{S^2}{n-1} V^T$$

$$\lambda_i = \frac{1}{n-1} S_i^2$$

$$\rightarrow \min_{W \in \mathbb{R}^{d,d}, U \in \mathbb{R}^{d,n}} \left\{ \sum_{i=1}^m \| \bar{x}_i - UW\bar{x}_i \|_2^2 \right\} = \underbrace{\sum_{i=n+1}^d \lambda_i(A)}_{\text{SUM}} \underbrace{\lambda_d, \lambda_{d-1}, \dots, \lambda_{n+1}, \lambda_n, \dots, \lambda_1}_{\text{Largest to Smallest}}$$

PCA PROCEDURE:

1. CENTER SAMPLES: $X \xrightarrow[n \times d]{ } X - \bar{\mu}$ / $\bar{\mu} = \frac{1}{m} \sum_{i=1}^m \bar{x}_i$ \rightarrow MEAN OF EACH FEATURE OF DATASET
2. APPLY PCA ON VECTORS $(\bar{x}_1 - \bar{\mu}), \dots, (\bar{x}_m - \bar{\mu})$

• PROPORTION OF VARIANCE EXPLAINED: PVE $\approx \frac{\lambda_i}{\sum_{j=1}^d \lambda_j} = \frac{\lambda_i}{\text{trace}(A)}$

• EFFICIENT IMPLEMENTATION OF PCA, FOR $d \gg m$:

$$\text{LET } B = X \cdot X^T \rightarrow B_{i,j} = \langle \bar{x}_i, \bar{x}_j \rangle$$

$$\cdot \text{IF } B\bar{u} = \lambda \bar{u} \rightarrow A(X^T \bar{u}) = X^T X X^T \bar{u} = X^T B \bar{u} = \lambda (X^T \bar{u})$$

$\rightarrow \frac{X^T \bar{u}}{\|X^T \bar{u}\|}$ IS AN EIGENVECTOR OF A, WITH EIGENVALUE λ

\rightarrow WE CAN PCA CALCULATING EIGENVALUES OF B INSTEAD OF A $\approx O(m^3 + m^2 d)$

m: # POINTS
d: # FEATURES

so, PCA $d = 4$

$$n = 5$$

f1	f2	f3	f4
1	2	3	4
5	5	6	7
1	4	2	3
5	3	2	1
8	1	2	2

$$\begin{aligned} \bar{M}_1 &= \frac{1+5+1+5+8}{5} = 4 \\ \bar{M}_2 &= \dots = 3 \\ \bar{M}_3 &= \dots = 3 \\ \bar{M}_4 &= \dots = 3.4 \end{aligned}$$

1. DATA ARE STANDARDIZED:

$$X \xrightarrow{} \frac{X - \bar{\mu}}{G} \rightarrow$$

\rightarrow WHERE ALSO G IS CONSIDERED

f1	f2	f3	f4
-1	-0.63246	0	0.26062
0.33333	1.26491	1.73205	1.56374
-1	0.63246	-0.57735	-0.17375
0.33333	0	-0.57735	-1.04249
1.33333	-1.26491	-0.57735	-0.60812

2. COMPUTE Σ :

$$\Sigma(f_i, f_j) = \frac{1}{n} \sum_{i=1}^n (\bar{f}_i - \bar{\mu}_i)(\bar{f}_j - \bar{\mu}_j)$$

$\hookrightarrow n=4$ FOR SAMPLES
 \hookrightarrow WHERE n USED

STANDARDIZED DATA SET

f1	f2	f3	f4
-1	-0.63246	0	0.26062
0.33333	1.26491	1.73205	1.56374
-1	0.63246	-0.57735	-0.17375
0.33333	0	-0.57735	-1.04249
1.33333	-1.26491	-0.57735	-0.60812

$$\rightarrow \text{Cov}(f_1, f_2) = \frac{1}{5} \left[(-1 - 0)(-0.63246 - 0) + \dots + (1.33 - 0)(-1.26491 - 0) \right] = -0.25298$$

• Σ :

	f1	f2	f3	f4
f1	0.8	-0.25298	0.03849	-0.14479
f2	-0.25298	0.8	0.51121	0.4945
f3	0.03849	0.51121	0.8	0.75236
f4	-0.14479	0.4945	0.75236	0.8

3. CALCULATE λ_i AND \tilde{V}_i :

EIGENVALUE EIGENVECTOR

$$A \tilde{v} = \lambda \tilde{v} \rightarrow (A - \lambda I) \tilde{v} = 0 \Leftrightarrow A - \lambda I = 0$$

$A = \Sigma$ • EIGENVALUES: $A - \lambda I = 0 \Leftrightarrow \det(A - \lambda I) = 0$

$$\hookrightarrow \text{SOL: } \lambda = [2.5757, 1.0653, 0.3439, 0.0250]$$

• EIGENVECTORS:

$$\tilde{v}_i : (A - \lambda_i I) \tilde{v}_i = 0 \rightarrow \begin{pmatrix} 0.800000 - \lambda & -0.25298 & 0.03849 & -0.14479 \\ -0.25298 & 0.800000 - \lambda & 0.511208 & 0.494498 \\ 0.03849 & 0.511208 & 0.800000 - \lambda & 0.752355 \\ -0.14479 & 0.494498 & 0.752355 & 0.800000 - \lambda \end{pmatrix} \times \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = 0$$

$$\hookrightarrow \text{SOL: } \begin{array}{cccc|ccc} 0.161960 & -0.917059 & -0.307071 & 0.196162 \\ -0.524048 & 0.206922 & -0.817319 & 0.120610 \\ -0.585896 & -0.320539 & 0.188250 & -0.720099 \\ -0.596547 & -0.115935 & 0.449733 & 0.654547 \end{array} \hookrightarrow V \text{ matrix}$$

$\tilde{v}_1 \quad \tilde{v}_2 \quad \tilde{v}_3 \quad \tilde{v}_4$

4. SORT EIGENVECTORS BY THEIR CORRESPONDING λ :

\rightarrow IN THIS CASE THEY ARE ALREADY SORTED

K=2 PRINCIPAL COMPONENTS

5. TRANSFORM DATA:

f1	f2	f3	f4
-1.000000	-0.632456	0.000000	0.260623
0.333333	1.264911	1.732051	1.563740
-1.000000	0.632456	-0.577350	-0.173749
0.333333	0.000000	-0.577350	-1.042493
1.333333	-1.264911	-0.577350	-0.608121

$$\begin{matrix} e1 & e2 & nf1 & nf2 \\ 0.161960 & -0.917059 & 0.014003 & 0.755975 \\ -0.524048 & 0.206922 & -2.556534 & -0.780432 \\ -0.585896 & -0.320539 & -0.051480 & 1.253135 \\ -0.596547 & -0.115935 & 1.014150 & 0.000239 \\ (4,2) & & 1.579861 & -1.228917 \\ (5,2) & & & \end{matrix}$$

$$\rightarrow \bar{x}_{i,\text{PCA}} = \bar{x}_i \cdot V[:, :k] \hookrightarrow \text{MULTIPLY } \bar{x}_i \text{ BY FIRST K EIGENVECTORS}$$

• RANDOM PROJECTIONS: $\rightsquigarrow O(n \log k) / k$: TARGET DIMENSIONALITY

IN SOME CASE, RECONSTRUCTION OF \bar{x} FROM $\tilde{x} = W\bar{x}$ IS NOT IMPORTANT,
BUT MAYBE IT IS IMPORTANT TO KEEP SOME PROPERTIES, SUCH AS:

① DO NOT DISTORT DISTANCES: $\|\bar{x}_i - \bar{x}_j\| \approx \|\tilde{x}_i - \tilde{x}_j\|, \forall i, j$

$$\rightarrow \frac{\|W\bar{x}_i - W\bar{x}_j\|}{\|\bar{x}_i - \bar{x}_j\|} \approx 1, \forall i, j$$

② $\forall \tilde{x} \in Q$, $Q = \{ \bar{x}_i - \bar{x}_j : i, j \in [m] \}$, WE HAVE $\frac{\|W\tilde{x}\|}{\|\tilde{x}\|} \approx 1$

RANDOM PROJECTION ↗

• $\bar{x} \xrightarrow{\text{d} \times 1} \xrightarrow{\text{d} \times n} \xrightarrow{n \times 1} W \bar{x}$ / W IS A RANDOM MATRIX: $W_{i,j} \sim N(0, \frac{1}{n})$

• LET \bar{w}_i BE THE i^{th} ROW OF W :

$$\rightarrow E[\|W\bar{x}\|^2] = \sum_{i=1}^n E[(\langle \bar{w}_i, \bar{x} \rangle)^2] = \sum_{i=1}^n \bar{x}^\top E[\bar{w}_i \bar{w}_i^\top] \bar{x} = \\ = n \bar{x}^\top \left(\frac{1}{n} I \right) \bar{x} = \|\bar{x}\|^2 \rightarrow \text{PROPERTY } ② \text{ SATISFIED}$$

\rightarrow IN FACT: $\|W\bar{x}\|^2 \sim \chi_n^2 \rightsquigarrow \text{CHI-SQUARED DISTRIBUT.}$

$$\rightarrow P \left[\left| \frac{\|W\bar{x}\|^2}{\|\bar{x}\|^2} - 1 \right| > \varepsilon \right] \leq 2e^{-\varepsilon^2 \frac{n}{6}} \quad \text{MOFFLING'S INEQUALITY}$$

• LEMMA (JOHNSON-LINDSTEDTSS LEMMA): \rightsquigarrow OBTAINED APPLYING THE UNION BOUND
OVER ALL VECTORS IN Q

LET Q BE A FINITE SET OF VECTORS IN \mathbb{R}^d

$$\text{LET } S \in (0, 1) \text{ AND } n \in \mathbb{Z} / \varepsilon = \sqrt{\frac{6 \cdot \log(2|Q|/S)}{n}} \leq 3 \quad n \rightsquigarrow \text{REDUCED DIMENSION}$$

PROBABILITIES

\rightarrow THEN $P \geq 1 - S$ OVER A CHOICE OF A RANDOM MATRIX $W \in \mathbb{R}^{n,d}$, $W_{i,j} \sim N(0, \frac{1}{n})$

$$, \text{WE HAVE: } \max_{\tilde{x} \in Q} \left| \frac{\|W\tilde{x}\|^2}{\|\tilde{x}\|^2} - 1 \right| < \varepsilon$$

$$\hookrightarrow \text{SETTING } \varepsilon = \sqrt{\frac{6 \cdot \log(2|Q|/S)}{n}} \leq 3 \rightarrow \text{PROPERTY } ② \text{ SATISFIED,}\\ \text{WITH C.I } 1 - S$$

COMPRESSED SENSING:

- PRIOR ASSUMPTION:

$$\begin{aligned} & \tilde{x} \approx U \bar{x} \quad / \quad U \text{ is ORTHONORMAL} \\ & \| \bar{x} \|_0 := |\{i : x_i \neq 0\}| \leq s, \quad s \ll d \end{aligned}$$

- HOW TO STORE \bar{x} :

- FIND $\tilde{\alpha} = U^T \bar{x}$, THEN SAVE NON-ZERO ELEMENTS OF $\tilde{\alpha}$
- REQUIRES ORDER OF $s \cdot \log(d)$ STORAGE
 \hookrightarrow CAN'T WE JUST MEASURE THE PART THAT WON'T BE THROWN AWAY?
 \rightarrow COMPRESSED SENSING

- COMPRESSED SENSING MAIN RESULTS:

- ANY SPARSE SIGNAL CAN BE FULLY RECONSTRUCTED IF IT WAS COMPRESSED BY $\tilde{x} \mapsto W \tilde{x}$ / W SATISFIES RESTRICTED ISOPERIMETRIC PROPERTY (RIP)
 \rightarrow IF A MATRIX HAS THIS PROPERTY \rightarrow IT HAS A LOW DISTORTION OF THE NORM OF ANY SPARSE REPRESENTABLE VECTOR
- RECONSTRUCTION CAN BE COMPUTED IN POLYNOMIAL TIME BY SOLVING A LINEAR PROGRAM
- A RANDOM $m \times d$ MATRIX IS LIKELY TO SATISFY RIP CONDITION IF m IS GREAT THEN THE ORDER OF $s \cdot \log(d)$

- RESTRICTED ISOPERIMETRIC PROPERTY:

A MATRIX $W \in \mathbb{R}^{m,d}$ IS (ε, s) -RIP IF $\forall \tilde{x} \neq 0 / \| \tilde{x} \|_0 \leq s$:

$$\rightarrow \left| \frac{\| W \tilde{x} \|_2^2}{\| \tilde{x} \|_2^2} - 1 \right| \leq \varepsilon \quad \Rightarrow \quad 1 - \varepsilon \leq \frac{\| W \tilde{x} \|_2^2}{\| \tilde{x} \|_2^2} \leq 1 + \varepsilon$$

- RIP MATRICES ALSO LOSSLESS COMPRESSION FOR SPARSE VECTORS:

- THEOREM:

Let $\Sigma < 1$, W be a $(\epsilon, 2\Sigma)$ -RIP matrix,
, \bar{x} be a vector / $\|\bar{x}\|_0 \leq s$, $\bar{v} = W\bar{x}$,
, $\tilde{x} \in \operatorname{argmin}_{\tilde{v}} \|\tilde{v}\|_0$:
 $\bar{v}^T W \tilde{x} = \bar{v}$
 $\rightarrow \tilde{x} = \bar{x}$ $\leadsto \bar{x}$ CAN BE PERFECTLY RECONSTRUCTED

- EFFICIENT RECONSTRUCTION:

IF WE FURTHER ASSUME THAT $\Sigma < \frac{1}{7 + \sqrt{2}}$:

↑
LINEAR PROGRAMMING
PROBLEM

$$\rightarrow \bar{x} = \underbrace{\operatorname{argmin}_{\tilde{v}} \|\tilde{v}\|_0}_{\tilde{v}: W\tilde{v} = \bar{v}} = \operatorname{argmin}_{\tilde{v}} \|\tilde{v}\|_1$$

• SUMMARY: WE CAN RECONSTRUCT AN SPARSE VECTOR EFFICIENTLY
BASED ON $O(s \log(d))$ MEASUREMENTS

- PCA VS RANDOM PROJECTIONS:

• RANDOM PROJECTIONS GUARANTEE PERFECT RECOVERY $\forall O\left(\frac{m}{\log(d)}\right)$ - SPARSE VECTORS

• PCA GUARANTEES RECOVERY IF ALL EXAMPLES ARE IN AN n -DIMENSIONAL SUBSPACE

\rightarrow

• IF DATA IS $\tilde{c}_1, \dots, \tilde{c}_n$: $\begin{cases} \text{RANDOM PROJECTIONS} \rightarrow \text{PERFECT} \\ \text{PCA} \rightarrow \text{FAIL} \end{cases}$

• ONLY SELECT FIRST n PC

• IF $d \gg n$, DATA EXACTLY IN n -DIM SUBSPACE: $\begin{cases} \text{RANDOM PROJECTIONS} \rightarrow \text{MIGHT FAIL} \\ \text{PCA} \rightarrow \text{PERFECT} \end{cases}$

\hookrightarrow e.g. DATA IN \mathbb{R}^2 AND TAKE $n=2$ PC

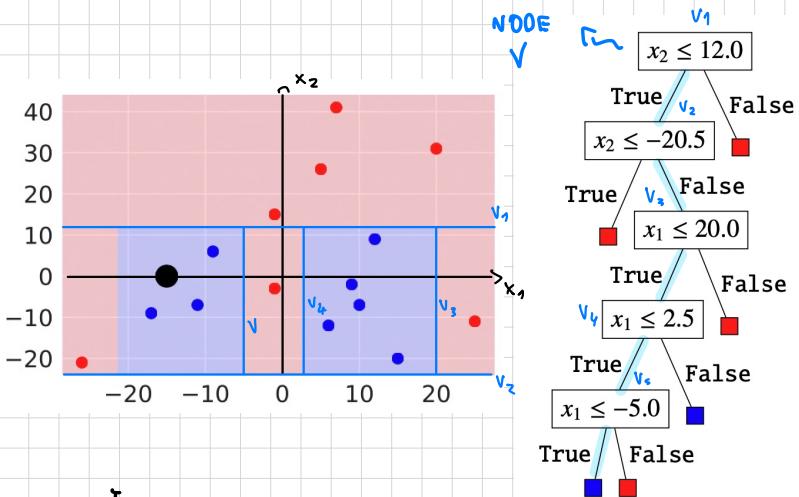
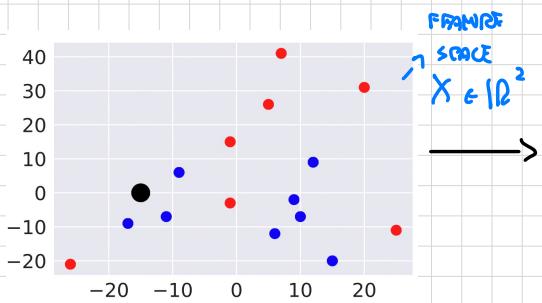
\rightarrow

IF LINEAR
DIMENSIONALITY:

$\tilde{x} \mapsto W\tilde{x}$

$\begin{cases} \text{PCA} : \text{OPTIMAL IF LINEAR RECONSTRUCTION AND ERROR IS SQUARED DISTANCE} \\ \text{REDUCTION;} \\ \text{RANDOM PROJECTIONS: PRESERVE DISTANCES, EXACT RECONSTRUCTION FOR} \\ \text{SPARSE VECTORS (BUT WITH NON-LINEAR RECONSTRUCTION)} \end{cases}$

DECISION TREES:



HOW TO CLASSIFY $\bullet = \bar{x} = [-15, 0]^T$:

→ DEFINE A FUNCTION $y(\bar{x}) : \bar{x} = [x_1, x_2] \mapsto \text{RED / BLUE}$

- ∀ NODE V , → REGION R_V OF FRAME SPACE X SPLITTED INTO 2 SUBREGIONS
- ∀ LEAF NODE W → REGIONAL PREDICTION FUNCTION y^W ON R_W
- ∀ TREE T → # REGIONS = $|W|$ / W : SET OF LEAF NODES

PREDICTION FUNCTION:

$$y(\bar{x}) = \sum_{w \in W} y^w(\bar{x}) \cdot \mathbb{1}_{[\bar{x} \in R_w]}$$

TRAINING LOSS:

GIVEN $y(\bar{x})$ AS ABOVE, TRAINING SET $\tilde{\gamma} = \{(\bar{x}_i, v_i)\}_{i=1}^n$:

$$l_{\tilde{\gamma}}(y) = \sum_{w \in W} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[\bar{x}_i \in R_w]} \cdot \underbrace{\text{Loss}(v_i, y^w(\bar{x}_i))}_{\text{CONTRIBUTION BY REGIONAL } y^w \text{ TO OVERALL TRAINING LOSS}}$$

CONSTRUCTION OF DECISION TREES:

IT IS NEEDED TO SPECIFY A SPLITTING RULE s FOR NODE v_i :

$$s : X \mapsto \{0, 1\}$$

↳ IT IS A CONDITION

ex. ROOT NODE OF PREVIOUS EXAMPLE: $\bar{x} \mapsto 1 | [x_1 \leq 12]$

USING A SPLITTING RULE s WE CAN DIVIDE AND

SUBSET $G \subseteq \tilde{\tau}$ OF THE TRAINING SET IN:

$$G_L \leftarrow G_T := \{(x, y) \in G : s(\bar{x}) = \text{TRUE}\}$$

$$G_R \leftarrow \{(x, y) \in G : s(\bar{x}) = \text{FALSE}\}$$

→ FINAL TREE CAN BE BUILT USING:

$$T = \text{CONSTRUCT_SUBTREE}(v_0, \tilde{\tau})$$

/ T_v : SUBTREE OF T FROM NODE v

$$/ y^v(\bar{x}) = y^L(\bar{x}) \mathbb{I}_{[\bar{x} \in R_{v_L}]} + y^R(\bar{x}) \mathbb{I}_{[\bar{x} \in R_{v_R}]}$$

Algorithm 1: Construct_Subtree
Input: A node v and a subset of the training data: $\sigma \subseteq \tau$.
Output: A (sub) decision tree T_v .

```

1 if termination criterion is met then      //  $v$  is a leaf node
2   Train a regional prediction function  $g^v$  using training data  $\sigma$ .
3 else
4   Find the best splitting rule  $s_v$  for node  $v$ .
5   Create successors  $v_L$  and  $v_R$  of  $v$ .
6    $\sigma_L \leftarrow \{(x, y) \in \sigma : s_v(x) = \text{True}\}$ 
7    $\sigma_R \leftarrow \{(x, y) \in \sigma : s_v(x) = \text{False}\}$ 
8    $T_{v_L} \leftarrow \text{Construct\_Subtree}(v_L, \sigma_L)$       // left branch
9    $T_{v_R} \leftarrow \text{Construct\_Subtree}(v_R, \sigma_R)$       // right branch
10 return  $T_v$ 
```

REGIONAL PREDICTION FUNCTION FOR CLASSIFICATION:

GIVEN LABELS $z = 0, \dots, c-1$, n_w : # FEATURE VECTORS IN R_w ,

LET $P_z^w = \frac{1}{n_w} \sum_{\{(x, y) \in \tilde{\tau} : x \in R_w\}} \mathbb{I}_{[y=z]}$ BE THE PROPORTION OF

FEATURE VECTORS IN R_w WITH LABEL $Z = 0, \dots, c-1$

→ REGIONAL PREDICTION FUNCTION FOR NODE w : $y^w(\bar{x}) = \underset{z \in \{0, \dots, c-1\}}{\text{argmax}} \{P_z^w\}$

REGIONAL PREDICTION FUNCTION FOR REGRESSION:

y^w IS USUALLY CHOSEN AS THE MEAN RESPONSE IN THAT REGION;

$$y^w(\bar{x}) = \bar{y}_{R_w} := \frac{1}{n_w} \sum_{\{(x, y) \in \tilde{\tau} : x \in R_w\}} y$$

- OPTIMAL SPLITTING RULES:

GIVEN $s(\bar{x}) = \mathbb{1}_{[x_j \leq \varepsilon]}$ WE WANT TO CHOOSE ε TO MINIMIZE:

$$\frac{1}{n} \sum_{j=1}^n \mathbb{1}_{[(\bar{x}_i, r) \in G_r]} \cdot \text{Loss}(y_i, y^r(\bar{x}_i)) + \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{[(\bar{x}_i, r) \in G_F]} \cdot \text{Loss}(y_i, y^F(\bar{x}_i))$$

- FOR REGRESSION, WE WANT TO MINIMIZE:

$$\frac{1}{n} \sum_{(\bar{x}, r) \in \tilde{\Gamma}: x_j \leq \varepsilon} (r - \bar{y}_r)^2 + \frac{1}{n} \sum_{(\bar{x}, r) \in \tilde{\Gamma}: x_j > \varepsilon} (r - \bar{y}_F)^2$$

LET $\{x_{j,k}\}_{k=1}^m$ DENOTE THE POSSIBLE VALUES OF x_j , $j = 1, \dots, p$
WITHIN THE TRAINING SUBSET $\tilde{\Gamma}$ ($m \ll n$ EXAMPLES)

→ TO MINIMIZE THE LOSS; EVALUATE IT $\forall m \times p$ VALUES $x_{j,k}$
→ THEN TAKE THE MINIMIZING PAIR $(j, x_{j,k})$

- FOR CLASSIFICATION, WE WANT TO MINIMIZE:

$$\frac{1}{n} \sum_{(\bar{x}, r) \in G_r} \mathbb{1}_{[y \neq y_r^*]} + \frac{1}{n} \sum_{(\bar{x}, r) \in G_F} \mathbb{1}_{[y \neq y_F^*]} \quad / \quad y^*: \text{most prevalent class (majority vote)}$$

- IMPURITIES:

MINIMIZING LOSS CAN ALSO SERVE AS MINIMIZING A WEIGHTED AVERAGE

OF "IMPURITIES" OF NODES G_r AND G_F :

$$\rightarrow \frac{1}{|G|} \sum_{(\bar{x}, r) \in G} \mathbb{1}_{[y \neq y^*]} = 1 - \frac{1}{|G|} \sum_{(\bar{x}, r) \in G} \mathbb{1}_{[r = y^*]} = 1 - p_r^* = \boxed{1 - \max_{z \in \{0, \dots, c-1\}} \{p_z\}}$$

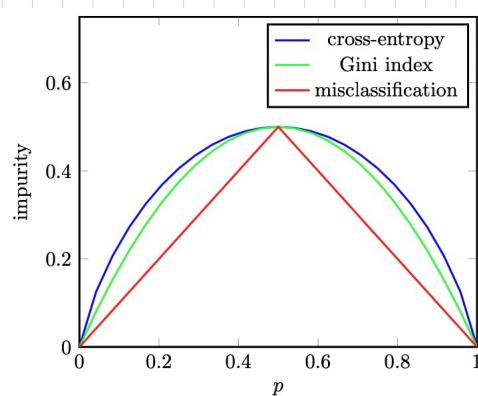
MISCLASSIFICATION IMPURITY

- OTHER IMPURITY MEASURES:

• ENTROPY : $- \sum_{z=0}^{c-1} p_z \cdot \log_2(p_z)$

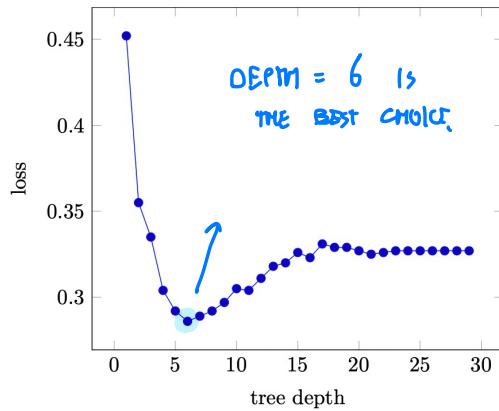
• GINI : $\frac{1}{2} \left(1 - \sum_{z=0}^{c-1} p_z^2 \right)$

→ MAX/MIN \Leftrightarrow LABEL PROPORTION IS $\frac{1}{c}$



- TERMINATION CRITERIA :

- # DATA POINTS IN TREE NODE \leq THRESHOLD
- MAXIMAL TREE DEPTH
- NO ADVANTAGE (OF TRAINING LOSS) IN SPLITTING REGIONS



ENSEMBLE METHODS:

BOOTSTRAP AGGREGATION:

GOAL: TRY TO REDUCE VARIANCE OF DTs TO ENHANCE DATA

- suppose B i.i.d. COPIES OF A TRAINING DATASET $\tilde{T} : \tilde{T}_1, \dots, \tilde{T}_B$

- WE CAN TAKE TRAIN B SEPARATE DTs, WITH AVERAGE PREDICTION:

$$y_{\text{avg}}(\bar{x}) = \frac{1}{B} \sum_{b=1}^B y_{\tilde{T}_b}(\bar{x}) \quad \text{EXPECTED TRAINER}$$

$$\rightarrow \text{BY LLN: } y_{\text{avg}}(\bar{x}) \xrightarrow[B \rightarrow \infty]{} y^E := E[y_{\tilde{T}}]$$

THEOREM:

EXPECTED SQUARE-ERROR GENERALIZATION RISK:

LET \tilde{T} BE A RANDOM TRAINING SET AND LET \bar{X}, Y BE A RANDOM FEATURE VECTOR AND RESPONSE THAT ARE INDEPENDENT OF \tilde{T} :

$$E(Y - y_{\tilde{T}}(\bar{X}))^2 \geq E(Y - y^E(\bar{X}))^2$$

- MULTIPLE INDEPENDENT DATASETS ARE RARELY AVAILABLE

\rightarrow WE CAN SUBSTITUTE THEM WITH BOOTSTRAPPED ONES:

WE CAN OBTAIN $\tilde{T}_1^*, \dots, \tilde{T}_B^*$ BY RESAMPLING FROM A SINGLE FIXED TRAINING SET \tilde{T} , AND USE THEM TO TRAIN B SEPARATE MODELS;

$\rightsquigarrow \approx \text{BOOTSTRAPPED}$

\rightarrow WE OBTAIN A BAGGED ESTIMATOR OF THE FORM:

$$y_{\text{bag}}(\bar{x}) = \frac{1}{B} \sum_{b=1}^B y_{\tilde{T}_b^*}(\bar{x})$$

LESS VARIANCE

• MOST EFFECTIVE FOR SENSITIVE MODELS

LITTLE CHANGE IN MEAN...

$$\cdot V(y_{\text{bag}}(\bar{x})) = \frac{1}{B} V(y_{\tilde{T}}(\bar{x}))$$

Algorithm 1: Bootstrap Aggregation Sampling

Input: Training set $\tau = \{(x_i, y_i)\}_{i=1}^n$ and resample size B .

Output: Bootstrapped data sets.

```

1 for  $b = 1$  to  $B$  do
2    $\mathcal{T}_b^* \leftarrow \emptyset$ 
3   for  $i = 1$  to  $n$  do
4     Draw  $U \sim \mathcal{U}(0, 1)$ ;
5      $I \leftarrow \lceil nU \rceil$  // select random index;
6      $\mathcal{T}_b^* \leftarrow \mathcal{T}_b^* \cup \{(x_I, y_I)\}$ .
7 return  $\mathcal{T}_b^*, b = 1, \dots, B$ .
```

- IT CAN BE SHOWN THAT, FOR LARGE SAMPLE SIZES,
ON AVERAGE ABOUT $\frac{1}{3} \left(e^{-1} \approx 0.37 \text{ ACCUR} \right)$ OF ORIGINAL SAMPLE POINTS
ARE NOT INCLUDED IN THE BOOTSTRAPPED SET $\tilde{\mathcal{I}}_b^*$, $b = 1, \dots, B$
- \rightarrow THESE SAMPLES ($\tilde{\mathcal{I}}_b^* \text{ OUT OF B SIZE}$) CAN BE USED FOR THE LOSS ESTIMATION
so. RUNNING CODE "BaggingExample.py"

DecisionTreeRegressor R² score = 0.575438224929718

Bagging R² score = 0.7612121189201985

Bagging OOB R² score = 0.7758253149069059

CORRELATED PREDICTIONS:

LET $\tilde{Z}_b = \tilde{y}_{\tilde{\mathcal{I}}_b^*}(\tilde{x})$, $b = 1, \dots, B$ BE i.i.d PREDICTION VALUES
FROM $\tilde{\mathcal{I}}_1^*, \dots, \tilde{\mathcal{I}}_B^*$

$$\rightarrow \text{SUPPOSE } V(Z_b) = \sigma^2 \rightarrow V(\tilde{Z}_b) = \frac{1}{B} \sigma^2$$

\rightarrow IF $\{\tilde{\mathcal{I}}_b^*\}$ ARE USED INSTEAD, $\{\tilde{Z}_b\}$ WILL BE CORRELATED

$\rightarrow \tilde{Z}_b$ ARE IDENTICAL DISTRIBUTED, BUT NOT INDEPENDENT

$$\rightarrow V(\tilde{Z}_b) = \underbrace{s \sigma^2}_{\text{CORRELATION CORP.}} + \underbrace{\sigma^2 \frac{(1-s)}{B}}_{B \rightarrow \infty}$$

RANDOM FORESTS :

BAGGING PROBLEM:

Suppose \exists features with good split \rightarrow it will be selected
At root level \rightarrow HIGHLY CORRELATED PREDICTIONS

$$\nabla \sum_{b=1}^B y_{T_b}^*$$

\hookrightarrow SOL: RANDOM FOREST \rightarrow ONLY A SUBSET OF FEATURES CONSIDERED DURING TREE CONSTRUCTION

$\rightarrow \nabla \sum_{b=1}^B \hat{y}_b^* \rightarrow DT$ BUILT WITH $m \ll p$ FEATURES FOR SPLITTING RULES

- DEFAULT VALUES FOR m : $m = \lceil \frac{1}{3}p \rceil$ OR $m = \lfloor \sqrt{p} \rfloor$

\hookrightarrow IN PRACTICE: m IS AN HYPERPARAMETER

\rightarrow WE LOSE INTERPRETABILITY

\hookrightarrow SOL: FEATURE IMPORTANCE

- FEATURE IMPORTANCE:

∇ INTERNAL NODE $v \rightarrow \Delta_{\text{loss}}(v)$ IN TRAINING LOSS $\xrightarrow{\text{DECREASE LOSS}}$

- FEATURE IMPORTANCE OF X_j :

IT TELLS HOW MUCH
FEATURE X_j HAS IMPACT ON
DECREASING THE LOSS
DURING BUILDING OF T

$$I_T(X_j) = \sum_{V \text{ INTERNAL } \in T} \Delta_{\text{loss}}(v) \cdot \mathbb{1}_{[X_j \text{ associated with } v]} , 1 \leq j \leq p$$

$$\rightarrow I_{RF}(X_j) = \frac{1}{B} \sum_{b=1}^B I_{T_b}(X_j) , 1 \leq j \leq p$$

• BOOSTING:

IT IMPROVES ACCURACY OF ANY LEARNING ALGORITHM

• BOOSTING FOR REGRESSION:

→ IN BOOSTING, PREDICTION FUNCTIONS ARE LEARNED SEQUENTIALLY:

1. START FROM A WEAK LEARNER y_0 FOR DATA $\tilde{\tau} = \{(\bar{x}_i, v_i)\}_{i=1}^n$

2. "BOOST" y_0 TO A NEW LEARNER: $y_1 := y_0 + h_1$

$$h_1 = \underset{h \in H}{\operatorname{arg\,min}} \left\{ \frac{1}{n} \sum_{i=1}^n \text{Loss}(v_i, y_0(\bar{x}_i) + h(\bar{x}_i)) \right\}$$

3. REPEAT FOR y_2, \dots, y_B

4. UNTIL THE BOOSTING PREDICTION FUNCTION: $y_B(\bar{x}) = y_0(\bar{x}) + \sum_{b=1}^B h_b(\bar{x})$

→ OR USING A STEP PARAMETER γ : $y_b := y_{b-1} + \gamma h_b$

• BOOSTING PROCEDURE WITH SQUARED ERROR LOSS:

• START WITH: $y_0 := \frac{1}{n} \sum_{i=1}^n v_i$

• $\forall b = 1, \dots, B$: $\tilde{\tau}_b := \{(\bar{x}_i, v_i - y_{b-1}(\bar{x}_i))\}$
RESIDUALS $e_i^{(b)}$

• FIT PREDICTION FUNCTION h_b ON $\tilde{\tau}_b$: $y_b(\bar{x}) = y_{b-1}(\bar{x}) + \gamma h_b(\bar{x})$

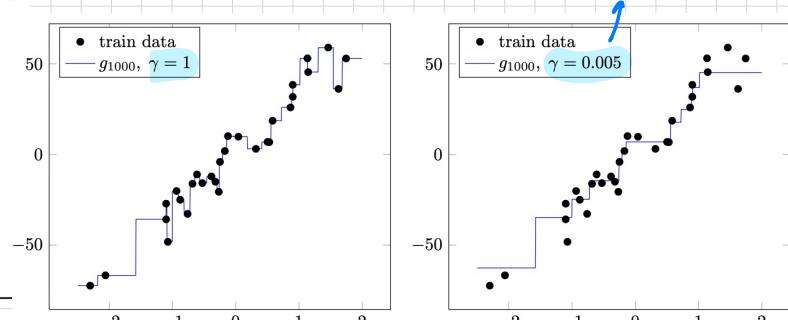
Algorithm 1: Regression Boosting with Squared-Error Loss

Input: Training set $\tau = \{(\bar{x}_i, v_i)\}_{i=1}^n$, the number of boosting rounds B , and a shrinkage step-size parameter γ .

Output: Boosted prediction function.

- 1 Set $g_0(\bar{x}) \leftarrow n^{-1} \sum_{i=1}^n v_i$.
- 2 **for** $b = 1$ **to** B **do**
- 3 Set $e_i^{(b)} \leftarrow v_i - g_{b-1}(\bar{x}_i)$ for $i = 1, \dots, n$, and let $\tau_b \leftarrow \{(\bar{x}_i, e_i^{(b)})\}_{i=1}^n$.
- 4 Fit a prediction function h_b on the training data τ_b .
- 5 Set $g_b(\bar{x}) \leftarrow g_{b-1}(\bar{x}) + \gamma h_b(\bar{x})$.
- 6 **return** g_B .

NO OVERFITTING



• GRADIENT BOOSTING :

- γ CAN BE VIEWED AS A STEP SIZE MADE IN THE OPPOSITE DIRECTION OF THE NEGATIVE GRADIENT OF SQUARED ERROR TRAINING LOSS :

$$z = \gamma_{b-1}(\bar{x}) \quad \begin{matrix} \uparrow \\ \rightarrow - \end{matrix} \quad \frac{\text{Loss}(y_i, z)}{\delta z} = \frac{\delta(y_i - z)^2}{\delta z} = \\ \vdots 2(y_i - \gamma_{b-1}(\bar{x}_i)) = 2e_i^{(b)}$$

- IT MIMICS A GRADIENT DESCENT ALGORITHM

Algorithm 2: Gradient Boosting

Input: Training set $\tau = \{(x_i, y_i)\}_{i=1}^n$, the number of boosting rounds B , a differentiable loss function $\text{Loss}(y, \hat{y})$, and a step-size parameter γ .

Output: Gradient boosted prediction function.

```

1 Set  $g_0(x) \leftarrow 0$ .
2 for  $b = 1$  to  $B$  do
3   for  $i = 1$  to  $n$  do
4     Evaluate the negative gradient of the loss at  $(x_i, y_i)$  via
       $r_i^{(b)} \leftarrow -\frac{\partial \text{Loss}(y_i, z)}{\partial z} \Big|_{z=g_{b-1}(x_i)} \quad i = 1, \dots, n$ .
5   Approximate the negative gradient by solving
       $h_b = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=0}^n (r_i^{(b)} - [g_{b-1}(x_i) + h(x_i)])^2$ .
6   Set  $g_b(x) \leftarrow g_{b-1}(x) + \gamma h_b(x)$ .
7 return  $g_B$ 
```

• ADA BOOST : \rightsquigarrow FOR BINARY $\{-1, 1\}$ CLASSIFICATION

- FINAL PREDICTION FUNCTION :

$$y_B(\bar{x}) = y_0(\bar{x}) + \sum_{b=1}^B h_b(\bar{x})$$

$$\quad / \quad h_b(\bar{x}) = \alpha_b c_b(\bar{x}), \quad \alpha_b \in \mathbb{R}_+, \quad c_b(\bar{x}) \in \{-1, 1\}$$

- \forall BOOSTING OPERATION :

$$(\alpha_b, c_b) = \underset{\alpha \geq 0, c \in C}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, y_{b-1}(\bar{x}_i) + \alpha c(\bar{x}_i)) \right\}$$

$$\quad / \quad \text{Loss}(y, \hat{y}) = e^{-y \cdot \hat{y}}, \quad \text{START} : y_0 := 0$$

(FURTHER ANALYSIS : $16 \div 19$)

- OPTIMAL α_b :

$$\alpha_b = \frac{1}{2} \ln \left(\frac{1 - \ell_\tau^{(b)}(c_b)}{\ell_\tau^{(b)}(c_b)} \right)$$

$$\quad / \quad \ell_\tau^{(b)}(c) := \frac{\sum_{i=1}^n w_i^{(b)} \mathbb{I}\{c(x_i) \neq y_i\}}{\sum_{i=1}^n w_i^{(b)}} \\ \quad \quad \quad \hookrightarrow \text{WEIGHTED 0-1 LOSS, ITERATION } b$$

- FINAL CLASSIFICATION :

$$\text{sign} \left(\sum_{b=1}^B \alpha_b c_b(x) \right)$$

Algorithm 3: AdaBoost

Input: Training set $\tau = \{(x_i, y_i)\}_{i=1}^n$, and the number of boosting rounds B .

Output: AdaBoost prediction function.

```

1 Set  $g_0(x) \leftarrow 0$ .
2 for  $i = 1$  to  $n$  do
3    $w_i^{(1)} \leftarrow 1/n$ 
4 for  $b = 1$  to  $B$  do
5   Fit a classifier  $c_b$  on the training set  $\tau$  by solving
       $c_b = \underset{c \in C}{\operatorname{argmin}} \ell_\tau^{(b)}(c) = \underset{c \in C}{\operatorname{argmin}} \frac{\sum_{i=1}^n w_i^{(b)} \mathbb{I}\{c(x_i) \neq y_i\}}{\sum_{i=1}^n w_i^{(b)}}$ .
6   Set  $\alpha_b \leftarrow \frac{1}{2} \ln \left( \frac{1 - \ell_\tau^{(b)}(c_b)}{\ell_\tau^{(b)}(c_b)} \right)$ . // Update weights
7   for  $i = 1$  to  $n$  do
8      $w_i^{(b+1)} \leftarrow w_i^{(b)} \exp\{-y_i \alpha_b c_b(x_i)\}$ .
9 return  $g_B(x) := \sum_{b=1}^B \alpha_b c_b(x)$ .
```