

INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 17

HASSAN ASHTIANI

THE VERSION WITH “BIAS”

Hard-SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

solve:

$$(\mathbf{w}_0, b_0) = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|^2 \text{ s.t. } \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \quad (15.2)$$

output: $\hat{\mathbf{w}} = \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}, \hat{b} = \frac{b_0}{\|\mathbf{w}_0\|}$

SOFT-MARGIN SVM

Soft-SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

parameter: $\lambda > 0$

solve:

$$\min_{\mathbf{w}, b, \xi} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right)$$

$$\text{s.t. } \forall i, y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

output: \mathbf{w}, b

$$\min_{\mathbf{w}} \left(\lambda \|\mathbf{w}\|^2 + L_S^{\text{hinge}}(\mathbf{w}) \right),$$

$$L_S^{\text{hinge}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y \langle \mathbf{w}, \mathbf{x}_i \rangle\}.$$

- WHAT IF WE NEED A NON-LINEAR CLASSIFIER?

SVM WITH BASIS FUNCTIONS

$$\phi(x) : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$$

- SOFT SVM IN OUR NOTATION

$$\min_{w \in \mathbb{R}^{d_1}} \frac{1}{n} \sum \max\{0, 1 - y^i \langle w, x^i \rangle\} + \lambda \|w\|_2^2$$

- SOFT SVM WITH BASIS FUNCTIONS

$$\min_{w \in \mathbb{R}^{d_2}} \frac{1}{n} \sum \max\{0, 1 - y^i \langle w, \phi(x^i) \rangle\} + \lambda \|w\|_2^2$$

- $\phi(x^i)$ (AND w) CAN BE HIGH-DIMENSIONAL
 - HOW TO DEAL WITH THE PROHIBITIVE COMPUTATIONAL COST?

REPRESENTOR THEOREM

- $W^* = \operatorname{argmin}_W \frac{1}{n} \sum \max\{0, 1 - y^i \langle w, \phi(x^i) \rangle\} + \lambda \|w\|_2^2$
- THEOREM: THERE ARE REAL-VALUES a_1, \dots, a_m SUCH THAT
$$W^* = \sum a_i \phi(x^i). \quad a^T \phi \quad \text{where } \phi = [\phi(x^1) \dots \phi(x^n)]^T$$
- INSTEAD OF d_2 PARAMS, CAN USE n PARAMS.
 - BETTER IF $d_2 \gg n$.

KERNELIZED SVM - PREDICTION

- $W^* = \operatorname{argmin}_W \frac{1}{n} \sum \max\{0, 1 - y^i \langle w, \phi(x^i) \rangle\} + \lambda \|w\|_2^2$

- $W^* = \sum a_i \phi(x^i)$.

- $K(x, z) = \langle \phi(x), \phi(z) \rangle$

efficiently computable

- $K(x) = (K(x, x^1), \dots, K(x, x^n))$.

$K_{n \times n} = \begin{bmatrix} \ddots & & \\ & \ddots & \\ & & K(x^1, x^1) \\ & & & \ddots \\ & & & & \ddots \end{bmatrix}$

- PREDICTION FOR A NEW TEST POINT USING \tilde{a} AND \tilde{K} :

$$\hat{y} = \operatorname{sgn}(\langle w^*, \phi(x) \rangle) = \operatorname{sgn}(\langle \sum a_i \phi(x^i), \phi(x) \rangle)$$

$$= \operatorname{sgn}(\sum a_i \underbrace{\langle \phi(x^i), \phi(x) \rangle}_{K(x, x^i)}) = \operatorname{sgn}(\tilde{a}^T \tilde{K}(x))$$

KERNELIZED SVM (LEARNING)

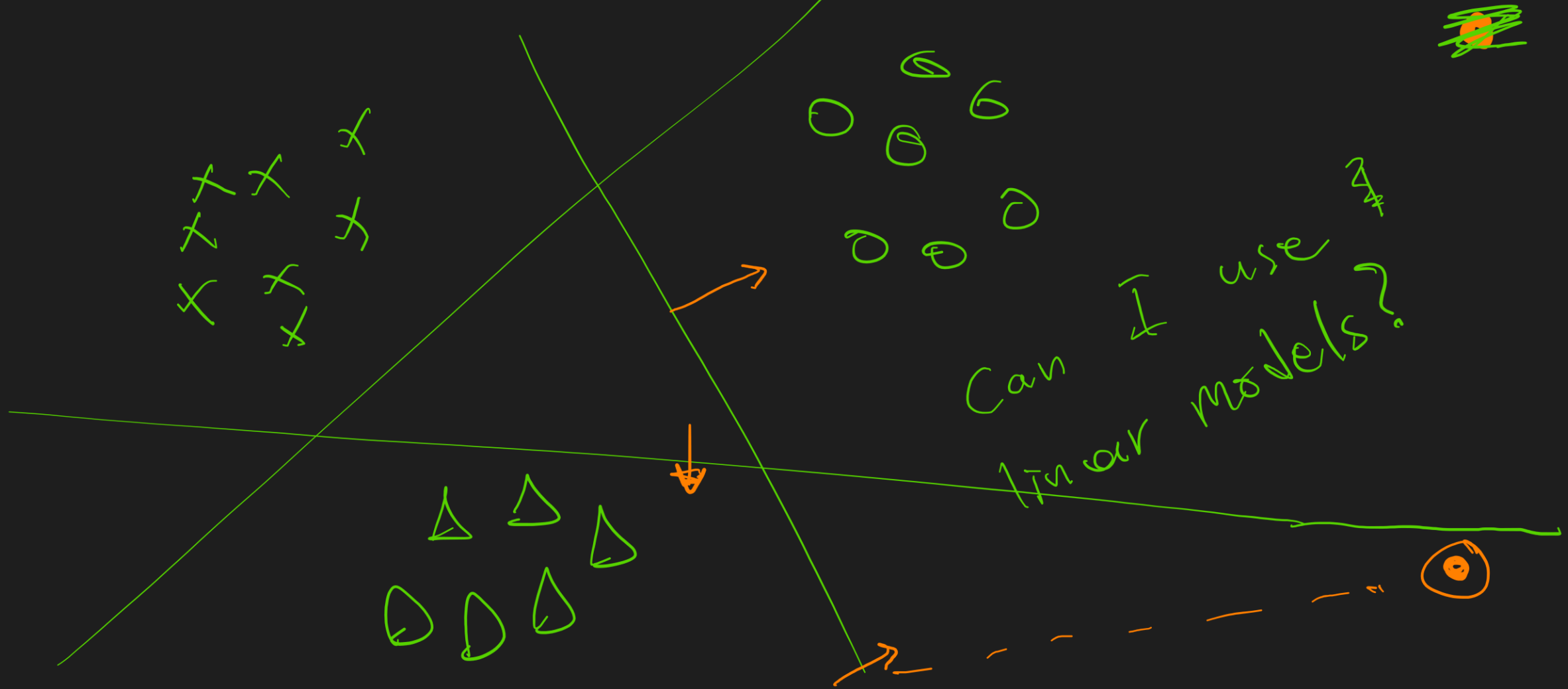
- $W^* = \underset{w}{\operatorname{argmin}} \frac{1}{n} \sum \max\{0, 1 - y^i \langle w, \phi(x^i) \rangle\} + \lambda \|w\|_2^2$

$$\operatorname{argmin}_a \left[\frac{1}{n} \sum_i \max\{0, 1 - y^i \langle \sum a_j \phi(x^j), \phi(x^i) \rangle\} + \lambda \|\sum a_i \phi(x^i)\|_2^2 \right]$$

$$= \operatorname{argmin}_{a \in \mathbb{R}^n} \left[\frac{1}{n} \sum_i \max\{0, 1 - y^i \langle a, K(x^i) \rangle\} + \lambda a^T K a \right]$$

$$\langle \sum a_i \phi(x^i), \sum a_j \phi(x^j) \rangle = \sum_{i,j} a_i a_j \langle \phi(x^i), \phi(x^j) \rangle = a^T K a$$

LINEAR MODELS FOR MULTICLASS CLASSIFICATION



ONE-VERSUS-ALL CLASSIFICATION

- TRAIN k DIFFERENT BINARY CLASSIFIERS

- CLASSIFIER i DISTINGUISHES SAMPLES FROM CLASS i VERSUS ALL OTHER CLASSES

- $h_i(x) = \text{SGN}(\langle w_i, x \rangle)$

- NOW FOR A NEW TEST POINT x

- $h^{\text{one-versus-all}}(x) = \text{argmax}_i(\langle w_i, x \rangle)$

binary
k classifiers

ALL-PAIRS CLASSIFICATION

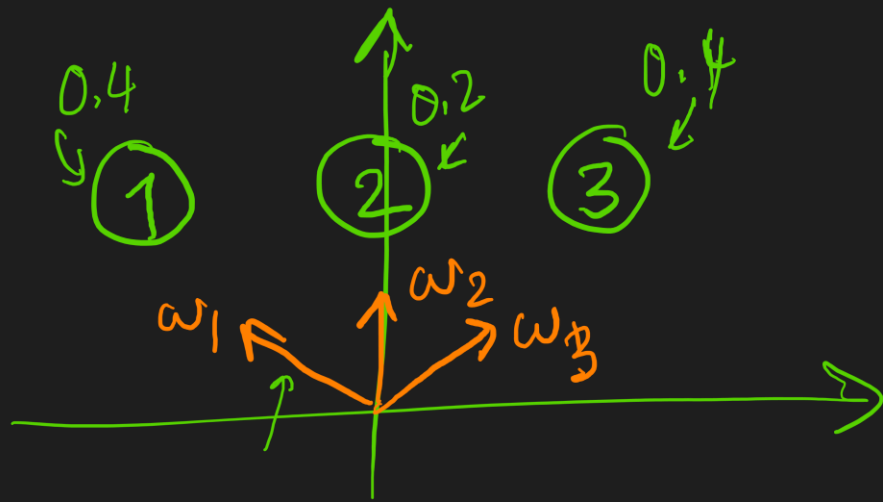
- FOR EACH DISTINCT $i, j \in \{1, 2, \dots, k\}$
 - TRAIN A CLASSIFIER TO DISTINGUISH SAMPLES FROM CLASS i AND SAMPLES FROM CLASS j
 - $h_{i,j}(x) = \text{SGN}(\langle w_{i,j}, x \rangle)$
- NOW FOR A NEW TEST POINT x
 - DO A VOTING AMONG $k(k - 1)/2$ CLASSIFIERS

if there was a tie, use the confidence values.

"GREEDY" VS "END-TO-END"

Show a scenario that the previous two approaches fail.

$$d=2, k=3$$



\mathbb{R}^2
 $[\omega_1 \quad \omega_2 \quad \omega_3] = W \in \mathbb{R}^6$
pick $\omega_1, \omega_2, \omega_3$ carefully
to get 0-error on
train data.

LINEAR MULTI-CLASS PREDICTOR?

$$x \in \mathbb{R}^d, y = 3$$

- THE MULTI-VECTOR ENCODING

→ • $y \in \{1, 2, \dots, k\}$

$$\Psi(x, y) = \left[\underbrace{0 \dots 0}_d \quad \underbrace{0 \dots 0}_d \quad \underbrace{x}_d \quad \underbrace{0 \dots 0}_k \right]$$

→ • (x, y) IS ENCODED AS $\Psi(x, y) = [0 \dots 0 \ x \ 0 \dots 0]^T \in \mathbb{R}^{d \times k}$

- $h(x) = \underset{y}{\operatorname{argmax}} \langle w, \Psi(x, y) \rangle$

↓
parameter

$$\langle w, \Psi(x, y) \rangle$$

END-TO-END VERSION OF ONE-VERSUS-ALL

Multiclass SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

parameters:

- regularization parameter $\lambda > 0$
- loss function $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$
- class-sensitive feature mapping $\Psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^{d \times K}$

solve:

e.g. $\mathcal{L}(\mathbf{w}) = \sum_{y \neq y_i}$

\mathbb{R}^d

≥ 0

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max_{y' \in \mathcal{Y}} (\Delta(y', y_i) + \langle \mathbf{w}, \Psi(\mathbf{x}_i, y') - \Psi(\mathbf{x}_i, y_i) \rangle) \right)$$

output the predictor $h_{\mathbf{w}}(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}, y) \rangle$

- (MORE IN CHAPTER 17 OF UNDERSTANDING MACHINE LEARNING)

CONFUSION MATRIX

A confusion matrix showing the relationship between predicted and actual classes. The matrix is annotated with blue circles around the values and blue arrows pointing to specific cells. The 'Actual class' header is at the top, and the 'Predicted class' header is on the left. The values in the matrix are: Cat (5, 2, 0), Dog (3, 3, 2), and Rabbit (0, 1, 11). The circles highlight the diagonal elements (5, 3, 0, 3, 1, 11) and the off-diagonal elements (2, 0, 2, 0, 1, 11). The arrows point to the 'Cat' column header, the 'Cat' row header, and the '0' value in the Cat-Rabbit cell.

		Actual class		
		Cat	Dog	Rabbit
Predicted class	Cat	5	2	0
	Dog	3	3	2
	Rabbit	0	1	11

wikipedia

- MORE DETAILED INFORMATION ABOUT WHICH CLASSES ARE BEING MISCLASSIFIED WITH WHICH

"RISK" MINIMIZATION

Multiclass SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

parameters:

regularization parameter $\lambda > 0$

loss function $\Delta : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

class-sensitive feature mapping $\Psi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$

solve:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max_{y' \in \mathcal{Y}} (\Delta(y', y_i) + \langle \mathbf{w}, \Psi(\mathbf{x}_i, y') - \Psi(\mathbf{x}_i, y_i) \rangle) \right)$$

output the predictor $h_{\mathbf{w}}(\mathbf{x}) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}, y) \rangle$

$$\Delta = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

		Actual class		
		Cat	Dog	Rabbit
Predicted class	Cat	5	2	0
	Dog	3	3	2
	Rabbit	0	1	11

$$\Delta = \begin{bmatrix} 0 & 0.1 & 0.5 \\ 0.1 & 0 & 0.8 \\ 0.6 & 0.7 & 0 \end{bmatrix}$$

DIFFERENT TYPES OF ERROR

- FALSE POSITIVE (TYPE I ERROR)
- FALSE NEGATIVE (TYPE II ERROR)

		Actual class	
		Cat	Non-cat
Predicted class	Cat	5 True Positives	2 False Positives
	Non-cat	3 False Negatives	17 True Negatives

you want to be sure that the phenomenon is present

- EXAMPLE: A CLASSIFICATION METHOD IS USED TO PREDICT WHETHER THERE IS A BRAIN TUMOR, BASED ON MRI DATA

- CASE 1: THE OUTCOME IS USED FOR DECIDING WHETHER TO DO A BRAIN SURGERY
- CASE 2: THE OUTCOME IS USED TO DECIDE WHETHER MORE TESTS (E.G., CT SCAN) IS REQUIRED

you don't want to miss

PRECISION VS RECALL

- TP = True Positive
- FP = False Positive
- TN = True Negative
- FN = False Negative

- $ACCURACY = \frac{TP+TN}{TP+TN+FP+FN}$
- $PRECISION = \frac{TP}{TP+FP}$, $RECALL = \frac{TP}{TP+FN}$
- ACCURACY IS NOT ALWAYS THE BEST MEASURE
- *BALANCED ERROR*
- $\frac{\alpha \cdot TP + (1-\alpha)TN}{TP+TN+FP+FN}$

Drawbacks of SVM?

- * Prediction-time complexity
- * need to store all training data
- * Dealing with $\tilde{K}_{n \times n}$ is difficult.
- * Choice of kernel is tricky
 - * often heuristic
 - * Not very data dependent