

COS234: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer



Topic: Ingredients of Machine Learning Problems

© 2020 YORAM SINGER

Bag of Words

- Need to classify an email as *spam* or *benign*
- Assume we are provided with a dictionary of “useful” words:
 $D = \{\text{clearance, order, singles, earn, money, cheap, cash, fast, dirt, diploma}\}$
- Index the words in some lexicographic order:
 $\text{clearance} \mapsto 1 \quad \text{order} \mapsto 2 \quad \dots \quad \text{diploma} \mapsto 10$
- Represent email as a bag-of-words vector:
 Entry j in b.o.w. vector (denoted \mathbf{x}) is **1** iff word $D[j]$ appears in email

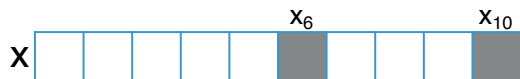
$$x[j] = 1 \Leftrightarrow D[j] \in \text{Email}$$

© 2020 YORAM SINGER

2

Bag of Words (example I)

“Dear dad, I would to get a scuba-diving diploma. The course ain’t cheap. Could I borrow \$1,000 until Aug?”



1 **0**

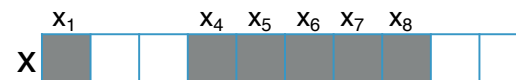
Note:
 $x[4] = x_4$

© 2020 YORAM SINGER

3

Bag of Words (example II)

“Dear, wouldn’t you like to earn some money fast? I got a box of clearance but genuine Rolexes which fell from the back of a driving truck. They are so cheap (\$9.99+tax) that I would get one for your mom & dad. Order quickly as these gems are selling fast. I accept only cash though.”



1 **0**

© 2020 YORAM SINGER

4

Spam/Benign Classification

- Suppose each word j is associated with a weight $w[j]$
- weight $w[j]$ represents relevance of j 'th word to "spaminess" of email:
 - $w[j]$ positive number \Rightarrow email containing $D[j]$ more likely to be **spam**
 - $w[j]$ negative number \Rightarrow email containing $D[j]$ more likely to be **benign**
- single word (typically) does not provide sufficient evidence
- compounding evidence and weighing "pros" and "cons"

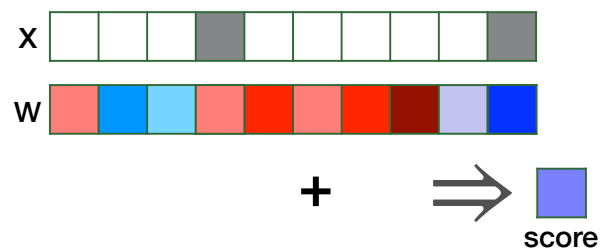
$$\sum_{j: D[j] \in \text{Email}} w[j] \Rightarrow \text{score}$$

© 2020 YORAM SINGER

5

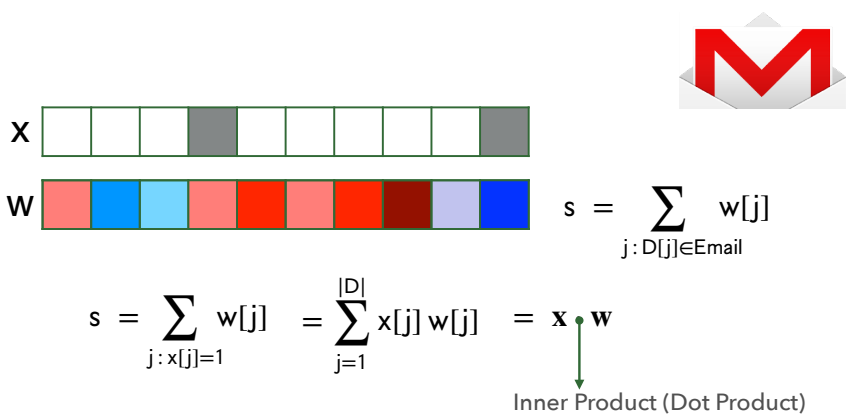
Back to Example 1

"Dear dad, I would to get a scuba-diving diploma. The course ain't cheap.
Could I borrow \$1,000 until Aug?"



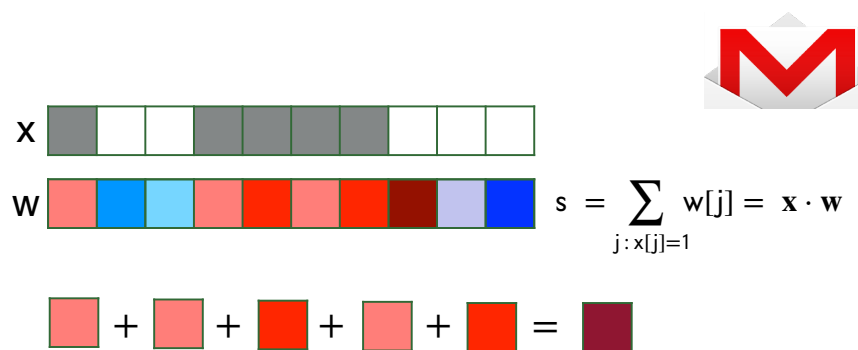
© 2020 YORAM SINGER

6



© 2020 YORAM SINGER

7

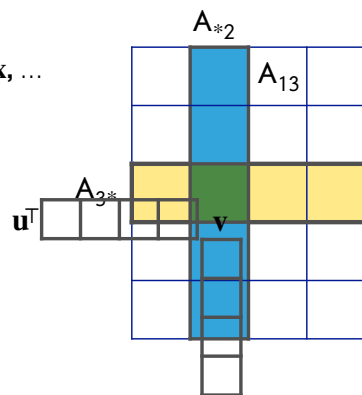


© 2020 YORAM SINGER

8

Note on Notation

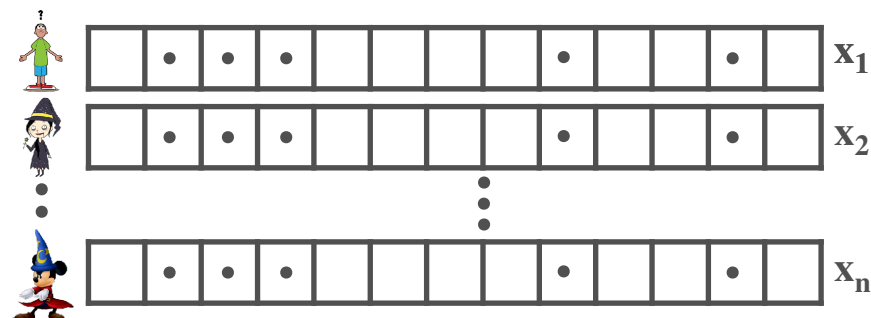
- [column] vectors in boldface: $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x}, \dots$
- i 'th element of vector \mathbf{v} : v_i
- element of array \mathbf{v} : $v[i]$ & $v['abc']$
- Matrices uppercase: A, W, X, \dots
- Element (i,j) of A : A_{ij}
- Row i of A : A_{i*} ; Column j of A : A_{*j}
- Inner (dot) product thus $\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v}$



© 2020 YORAM SINGER

9

Training Data



© 2020 YORAM SINGER

10

Training Data as MATRIX

#examples $n=5$

#features (dimension) $d=4$

X_{1*}				
X_{2*}				
X_{3*}				
X_{4*}				
X_{5*}				

Classification

+	y_1
+	y_2
-	y_3
+	y_4
-	y_5

labels

Regression

0.1	y_1
-3	y_2
-4	y_3
7.3	y_4
-12	y_5

targets

© 2020 YORAM SINGER

11

Test Data: MATRIX With Unknown Outputs

#test examples $n=5$

#features (dimension) $d=4$

X_{1*}^{te}				
X_{2*}^{te}				
X_{3*}^{te}				
X_{4*}^{te}				

Classification

?	y_1
?	y_2
?	y_3
?	y_4

labels

Regression

?	y_1
?	y_2
?	y_3
?	y_4

targets

© 2020 YORAM SINGER

12

Goal of Learning (linear case)

- Given data matrix \mathbf{X} and target vector \mathbf{y}

- Find weight vector \mathbf{w} such that as often as possible ?

$$\mathbf{w} \cdot \mathbf{x}_i \approx y_i$$

- Goodness of fit between predicted outcome and observed for all data ?

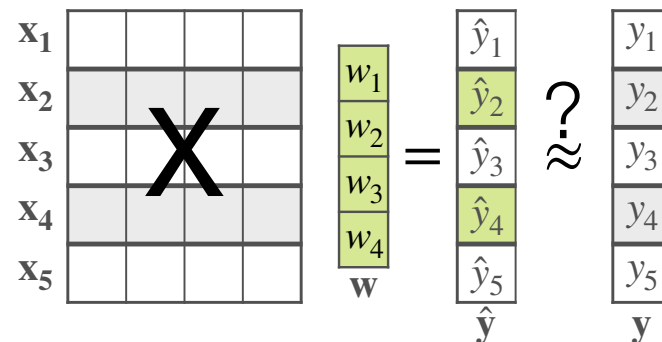
$$\mathbf{X}\mathbf{w} \approx \mathbf{y} \quad ?$$

- And, importantly, also want good predictions on unseen (test) data \mathbf{X}^{test}

© 2020 YORAM SINGER

13

Re-enter MATRIX



© 2020 YORAM SINGER

14

≈

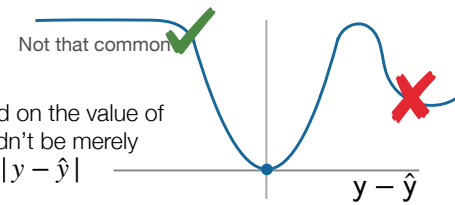
- Nobody is perfect:
 - input-output mapping is typically non-linear
 - noise in inputs & outputs
 - inherent ambiguity [in class experiment]
- Need to assess goodness of fit between y and \hat{y}
- Use a loss function: $\ell : \mathcal{R} \times \mathcal{R} \rightarrow \mathcal{R}_+$
- For now, classification aside and focus on regression with $y, \hat{y} \in \mathcal{R}$

© 2020 YORAM SINGER

15

Properties of ℓ

- When $y = \hat{y}$ loss should be 0
- When $y \neq \hat{y}$ loss should be ≥ 0
- If $|y_2 - \hat{y}_2| > |y_1 - \hat{y}_1|$ we typically want $\ell(y_2, \hat{y}_2) > \ell(y_1, \hat{y}_1)$



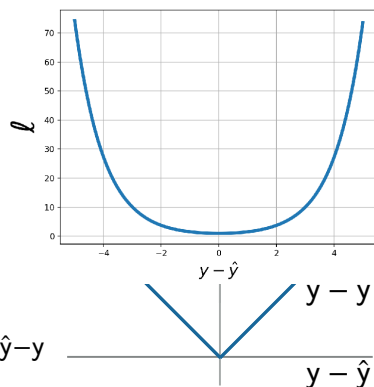
- May depend on the value of y and wouldn't be merely function of $|y - \hat{y}|$

© 2020 YORAM SINGER

16

Regression Losses

- Squared loss $\ell(y, \hat{y}) = (y - \hat{y})^2$
- Absolute loss $\ell(y, \hat{y}) = |y - \hat{y}|$
- Exponential loss $\ell(y, \hat{y}) = e^{y - \hat{y}} + e^{\hat{y} - y}$



© 2020 YORAM SINGER

17

Symmetrization of Losses

- Given $f : \mathfrak{R} \rightarrow \mathfrak{R}$ bounded from below: $\exists c : f(z) > c$

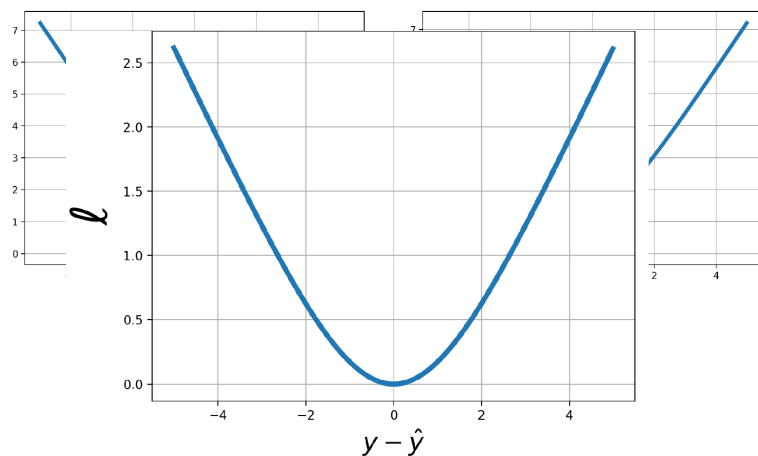
- Symmetrization at 0:

$$\frac{1}{2}(f(z) + f(-z)) - f(0)$$

- Use $z = y - \hat{y}$
- Exp-Loss: $f(z) = e^z$
- Log-Loss: $f(z) = \log(1 + e^z)$

© 2020 YORAM SINGER

18



© 2020 YORAM SINGER

19

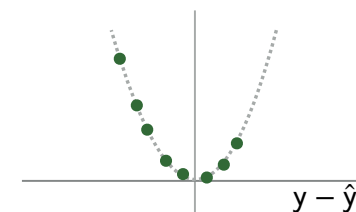
Training Loss

Average (why average?) loss over all examples

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i(\mathbf{w}))$$

For squared loss we can write

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$



© 2020 YORAM SINGER

20