

## Supervised Learning for Binary Classification

### - Money

Pos: banknote, penny  $\rightarrow$  dime is pos  
 Neg: bottle, pen

### - Movie preferences

Pos: thumbs up  $\rightarrow$  movie's T.J. likes  
 Neg: thumbs down

How would you classify  $x' = (\text{complete}, \text{no}, \text{no}, \text{clear}, \text{no})$

$\rightarrow$  inductive bias for today: find closest neighbors in training data

$\rightarrow x'$  is similar to  $x_2 \rightarrow$  YES

Notation:

$$1[P] = \begin{cases} 1 & \text{if } P \text{ is true} \\ 0 & \text{if } P \text{ is false} \end{cases}$$

$\operatorname{argmax}_{y \in Y} \{g(y)\}$  is the value  $y \in Y$  that maximizes  $g(y)$

Similarity measure:  $k(x, x') = \text{number of matching features}$

Example 1:  $x' = (\text{complete}, \text{no}, \text{no}, \text{clear}, \text{no})$

1 NN( $x'$ ) =  $\{2\}$  with  $k(x', x_2) = 4 \rightarrow$  vote  $y_2 = \text{yes} \rightarrow$  classify  $y' = \text{yes}$

2 NN( $x'$ ) =  $\{2, 1\}$  with  $k(x', x_1) = 3 \rightarrow$  two votes for yes  $\rightarrow y' = \text{yes}$

Example 2:  $x' = (\text{partial}, \text{yes}, \text{no}, \text{clear}, \text{no})$

2 NN( $x'$ ) =  $\{3, 1\} \rightarrow$  1 vote yes and 1 vote no  $\rightarrow y' = \text{flip coin}$

Similarity weighted kNN:

Example:  $x' = (\text{complete}, \text{yes}, \text{no}, \text{clear}, \text{no})$

Most similar:  $x_1 \rightarrow k(x_1, x') = 1$   
Second similar:  $x_2 \rightarrow$

---

KNN for Real-valued attributes

- Gaussian:  $k(x_i, x') = e^{-(x_i - x')^2}$

$\rightarrow k(x_i, x') = 1$   $x_i = x'$

$\rightarrow k(x_i, x') \rightarrow 0$  as  $x_i$  and  $x'$  increase in Euclidean dist

- Cosine:  $k(x_i, x') = \cos(x_i, x') + 1$

---

Supervised Learning Naming

if  $y \in \{-1, +1\}$ , then called binary classification

if  $y$  is discrete and size greater 2, then multi-class classification

if  $y$  is real number, then called regression

if  $y$  is combinatorial object (eg tree), structured output prediction

---

KNN Advantages and Disadvantages

simple

flexible

computationally expensive

keep all data