

# k-Nearest Neighbor

#### Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Eric Eaton (UPenn), David Kauchak (Pomona), and Andrew Moore (CMU), and the many others who made their course materials freely available online.

### **Data Representation**

Learning Goals

- Describe how to represent complex data
- View data graphically







- Well, we could map Sunny=0, Overcast=1, Rainy=2...
- But such a mapping may not always be appropriate
  - imagine feature values being red, blue, green
  - red=0, blue=1, green=2 implies red more similar to blue than to green
- Solution: for feature with K > 2 possible values, create binary features, one for each possible value

	<b>Y</b>	Out	т	R		Υ	(S?, O?, R?, T, R )
F	С	Sunny	Low	Yes		1	$\langle 1 , 0 , 0 , 0 , 1 \rangle$
1	N	Sunny	High	Yes		0	(1,0,0,1,1)
1	۱I	Sunny	High	No		0	(1,0,0,1,0)
F	>	Overcast	Low	Yes	$\Rightarrow$	1	<pre>(0, 1, 0, 0, 1)</pre>
F	>	Overcast	High	No		1	(0,1,0,1,0)
F	<b>&gt;</b>	Overcast	Low	No		1	<pre>(0, 1, 0, 0, 0)</pre>
1	N	Rainy	Low	Yes		0	(0,0,1,0,1)
F	>	Rainy	Low	No		1	<pre>(0,0,1,0,0)</pre>











- Binary-valued features
  - Hamming distance  $dist(\mathbf{a}, \mathbf{b}) = \sum_i I(a_i \neq b_i)$ counts number of features where two examples disagree
- Mixed feature types (some real, some binary)
  - mixed distance measures
  - e.g. Euclidean for real part, Hamming for binary part
- Can also assign weights to features  $dist(a, b) = \sum_i w_i \cdot d(a_i, b_i)$



• Can be changed by different distance metrics





 $dist(a,b) = (a_1 - b_1)^2 + (a_2 - b_2)^2$ 

 $dist(a,b) = (a_1 - b_1)^2 + (3a_2 - 3b_2)^2$ 

• Become more complex as more examples are stored

Based on slide by Eric Eaton (originally by Andrew Moore)





# *k*-Nearest Neighbor (*k*-NN)

To classify example **x**:

- Find k nearest neighbors of x
- Choose as label the majority label within k nearest neighbors

#### How do we choose k?

- Often data-dependent and heuristic-based
  - common heuristic : choose 3,5,7 (odd number)
- Use validation data
- In general, k too small or too big is bad

Based on slide by David Kauchak and Piyush Rai

# k-Nearest Neighbor (k-NN)

To classify example **x**:

- Find k nearest neighbors of x
- Choose as label the majority label within *k* nearest neighbors

#### Any variants?

- Fixed distance
  - instead of k-NN, count majority from all examples within fixed distance (radius-based neighbors)
- Weighted
  - instead of treating all examples equally, weight "vote" of examples so that closer examples have more vote/weight (often use some sort of exponential decay)

### kNN Problems and ML Terminology

Learning Goals

- Describe how to speed-up kNN
- Define non-parametric and parametric and describe differences
- Describe curse of dimensionality

## Speeding up *k*-NN

- *k*-NN is a "lazy" learning algorithm
  does virtually nothing at training time
- But classification / prediction can be costly when training set is large
  - for *n* training examples and *d* features, how many computations required for each test example?
- Two strategies for alleviating this weakness
  - edited nearest neighbor : do not retain every training instance
  - k-d tree : use smart data structure to look up NN







### Summary: k-NN

#### When to consider

- examples map to points in  $\mathbb{R}^d$
- small number (<20) attributes per instance
- lots of training data

#### **Advantages**

- "training" is very fast
- simple to implement
- learn complex target functions
- adapts well to online learning
- do not lose information
- robust to noisy training data (when k > 1)

#### **Disadvantages**

- slow prediction
- easily fooled by irrelevant attributes
- sensitive to range of feature values
- not much insight into problem domain because no explicit model

Based on slides by David Sontag and David Page