

The nature of data unit:

- ▶ Objects and vectors (Wednesday)
- ▶ K nearest neighbors classification (**today**)
- ▶ (Start linear regression next week)

Today:

- ▶ Leftover: Meet the data; features as functions
- ▶ Concept
- ▶ Algorithm and analysis
- ▶ Things to notice
- ▶ Distance metrics and norms

*Let  $X$  be the training data of  $N$   $D$ -dimensional observations.*

*Let  $k$  be a natural number.*

*Given a new datapoint  $\vec{x}$ , assign  $\vec{x}$  to the same class as the majority of  $k$  training datapoints that are closest to  $\vec{x}$  in the vector space.*

*Instance-based methods can also use more complex, symbolic representations for instances. . . . Case-based reasoning has been applied to tasks such as storing and reusing past experience at a help desk, reasoning about about legal cases by referring to previous cases, and solving complex scheduling problems by reusing relevant portions of previously solved problems.*

*Mitchell, Machine Learning, p 231.*

A *metric* or *distance function* between two vectors is a function  $d : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  with the properties that for any vectors  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^D$ ,

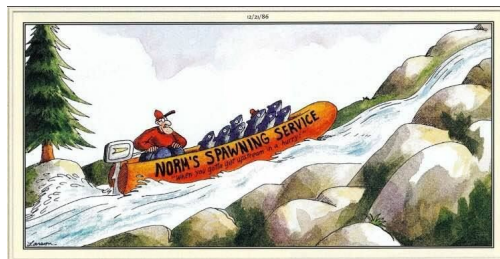
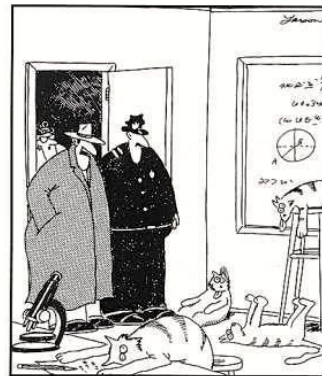
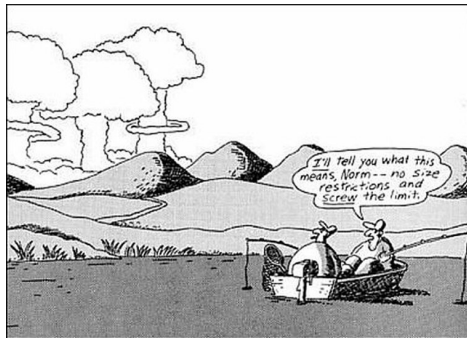
- ▶  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$  (symmetry)
- ▶  $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$  (the triangle inequality)
- ▶  $d(\mathbf{x}, \mathbf{y}) = 0$  iff  $\mathbf{x} = \mathbf{y}$  (the identity of indiscernibles)

A *norm* is a function  $\|\cdot\| : \mathbb{R}^D \rightarrow \mathbb{R}$  with the properties that for any vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$  and any scalar  $\lambda \in \mathbb{R}$ ,

- ▶  $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\|$  (absolute homogeneity)
- ▶  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$  (the triangle inequality)
- ▶  $\|\mathbf{x}\| \geq 0$  and, moreover,  $\|\mathbf{x}\| = 0$  iff  $\mathbf{x} = \mathbf{0}$  (positive definiteness)

Any norm induces a metric with

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|$$



Euclidean distance ( $L_2$  norm):

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(x_0 - y_0)^2 + (x_1 - y_1)^2}$$

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\left(\sum_{i=0}^{D-1} (x_i - y_i)^2\right)}$$

Manhattan or city-block distance ( $L_1$  norm):

$$d(\mathbf{x}, \mathbf{y}) = |x_0 - y_0| + |x_1 - y_1|$$

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{D-1} |x_i - y_i|$$

Minkowski distance ( $L_p$  norm):

$$d(\mathbf{x}, \mathbf{y}) = (|x_0 - y_0|^p + |x_1 - y_1|^p)^{\frac{1}{p}}$$

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=0}^{D-1} |x_i - y_i|^p\right)^{\frac{1}{p}}$$

*Mahalanobis distance:* Let  $\mathbf{S}$  be the covariance matrix of the entire dataset  $\mathbf{X}$ , and so  $\mathbf{S}^{-1}$  is the inverse of the covariance matrix.

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{S}^{-1} (\mathbf{x} - \mathbf{y})}$$

Compare with Euclidean distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\left( \sum_{i=0}^{D-1} (x_i - y_i)^2 \right)} = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})}$$

*Canberra distance:*

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{D-1} \frac{|x_i - y_i|}{|x_i| + |y_i|}$$

Compare with Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = |x_0 - y_0| + |x_1 - y_1| \qquad d(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^{D-1} |x_i - y_i|$$



## Coming up:

### **Due Fri, Jan 24:**

*Do numpy programming assignment*

*Take the KNN quiz*

### **Due Wed, Jan 29:**

*Read and respond to article about the Boston Housing Dataset*

*(Find pdf on Canvas.)*

### **Due Thurs, Jan 30:**

*Read the textbook from Chapters 1 and 5 (see Canvas for specific sections)*

### **Due Fri, Jan 31:**

*Do KNN programming assignment*

### **Due Mon, Feb 3:**

*Propose project topic*