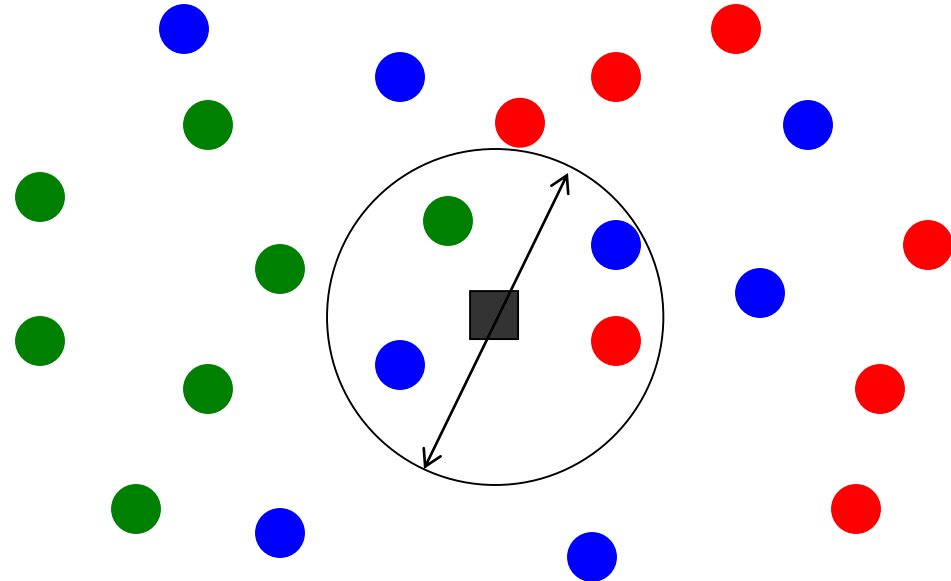


# The Nearest-Neighbor Classifier

- Let  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  be a set of “labeled” *training* samples (or *prototypes*)
- This means that each  $\mathbf{x}_i \in \omega_j$  (a class)
- What is the rule?... quite simple:
  - Given an **unknown** sample  $\mathbf{x}$
  - Decide the class of sample  $\mathbf{x}'$  that is **closest** to  $\mathbf{x}$
  - This is known as the 1-NN rule
- What is the definition of “closeness” ?
- For this, we need a “metric” to measure the “distance” between two samples (later)

# The $k$ -NN rule

- Start with sample  $\mathbf{x}$  we want to classify
- Grow a region until it encloses  $k$  samples  
or the  $k$  nearest neighbors
- Obtaining:
  - $k_1$  samples  $\in \omega_1$
  - $k_2$  samples  $\in \omega_2$
  - $\vdots$
  - $k_c$  samples  $\in \omega_c$
- Decide  $\omega_j$ 
  - where  $k_j$  is the max. of  
 $k_1, k_2, \dots, k_c$



Example: 2 classes –  $k$  must be odd to avoid ties

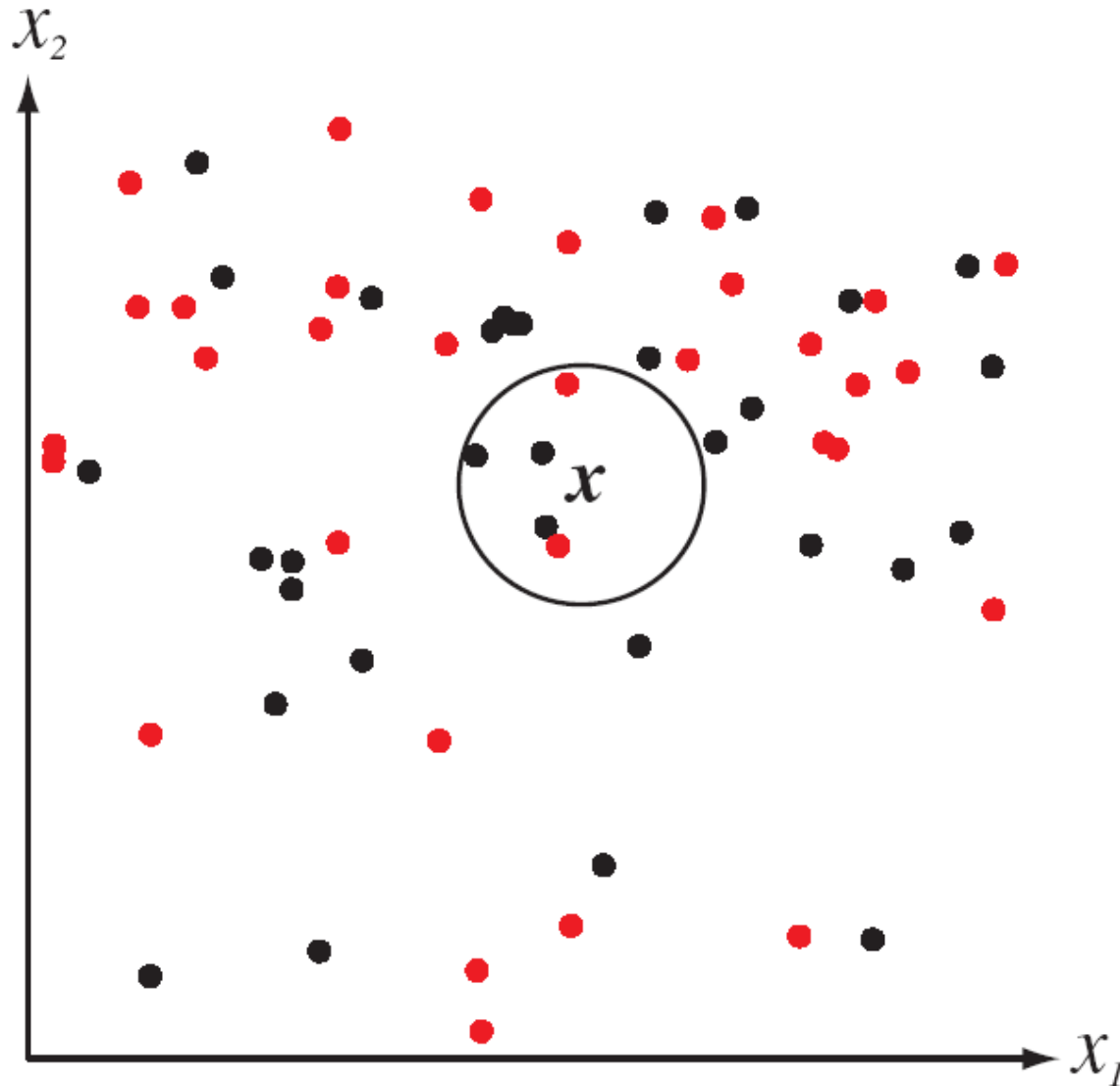
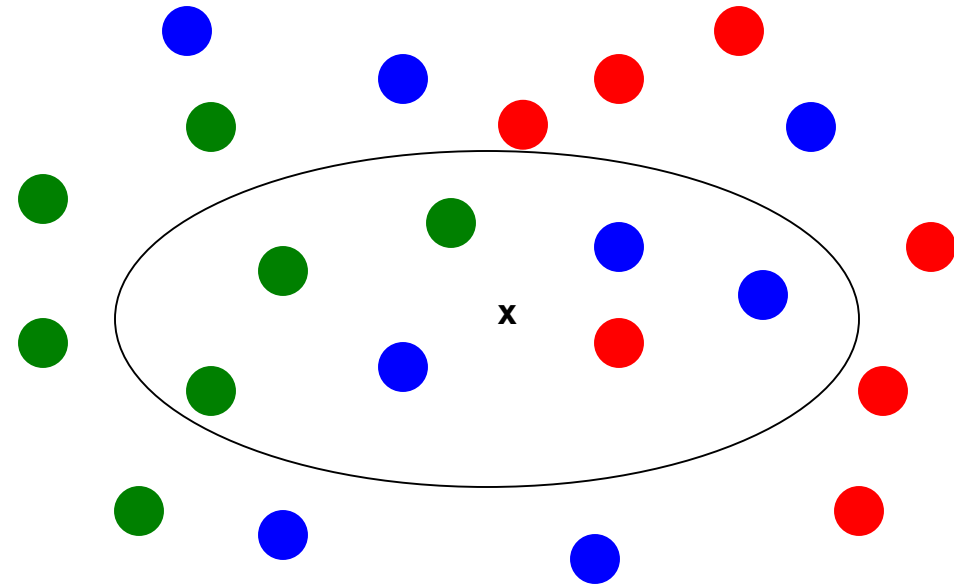


Figure from Duda et al.

## Example: Multi-class – ties resolved *arbitrarily*

- Distance: Mahalanobis
- Classify: **x**
- $k = 7$ 
  - $k_1 = 3$
  - $k_2 = 1$
  - $k_3 = 3$



- ... results in a tie between  $\omega_1$  and  $\omega_3$
- ... design a rule that resolves the ties!

# Proximity Measures - Metrics

- Proximity can be seen as a generalization of:
  - Dissimilarity (distance)
  - Similarity (closeness)
- Similarity and dissimilarity functions can satisfy certain properties.
  - These are called “metrics” if the following properties are satisfied...
- **Properties of metrics (distance):**
  - Let  $\mathbf{x}, \mathbf{y}, \mathbf{z}$  be vectors in  $\mathbb{R}^d$
  - Let  $d(\mathbf{x}, \mathbf{y})$  be the distance between  $\mathbf{x}$  and  $\mathbf{y}$
- **Reflexivity:**
$$d(\mathbf{x}, \mathbf{y}) = 0 \quad \text{iff} \quad \mathbf{x} = \mathbf{y}$$
- **Symmetry:**
$$d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$$
- **Positivity (nonnegativity):**
$$d(\mathbf{x}, \mathbf{y}) \geq 0$$
- **Triangle inequality:**
$$d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$$

# Distance functions (not all are metrics)

- **Euclidean distance:**

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

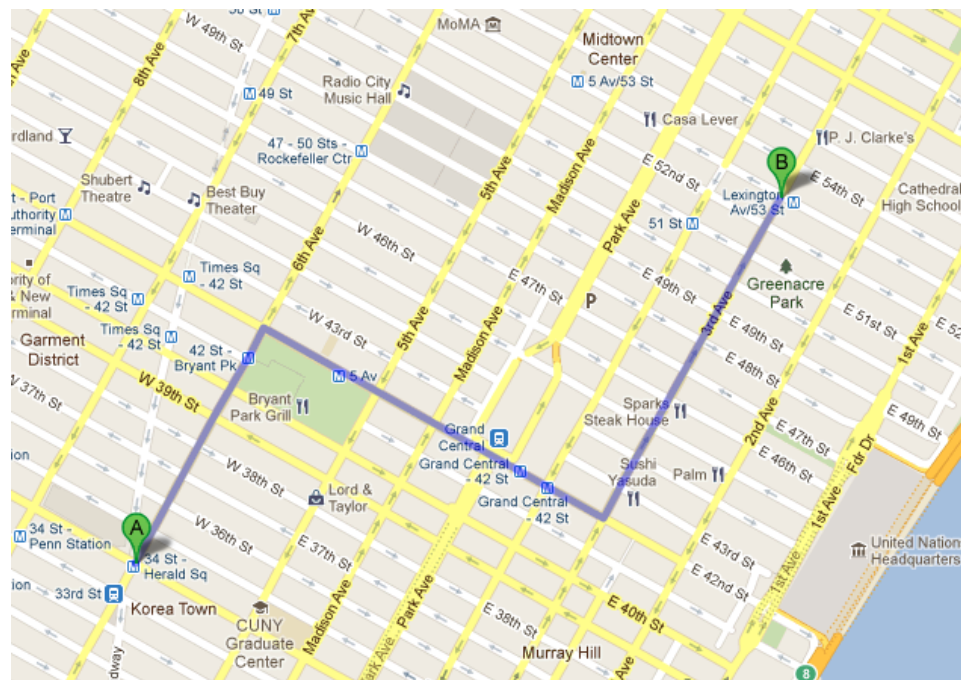
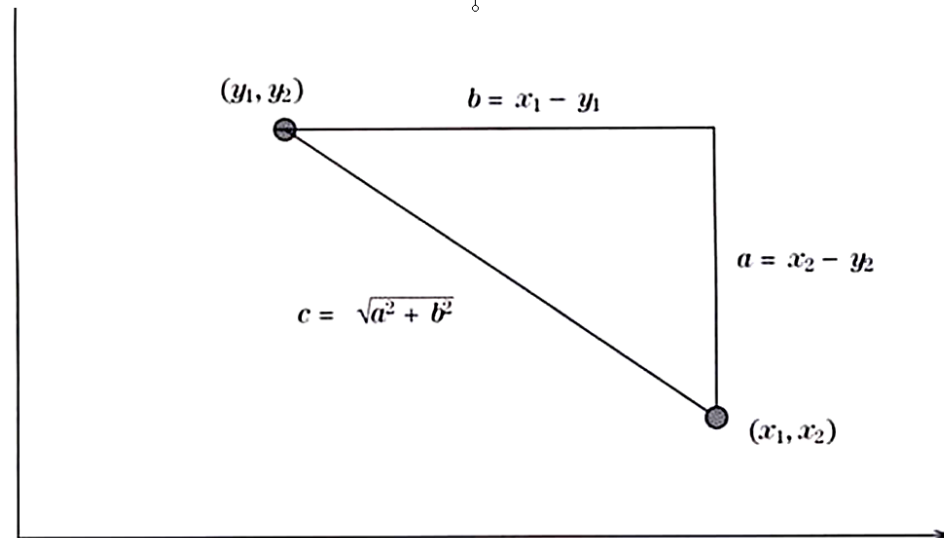
- **Minkowski distance:**

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

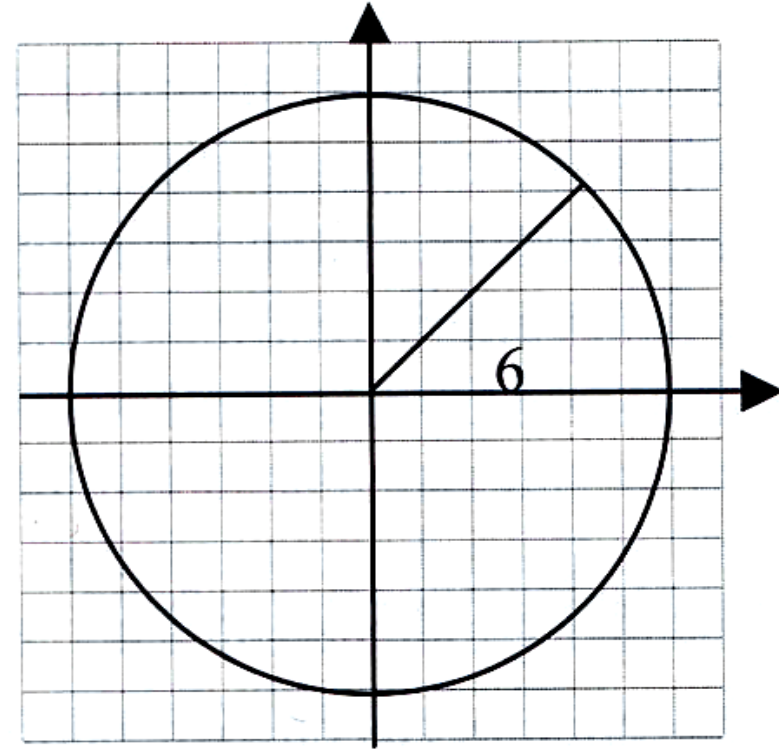
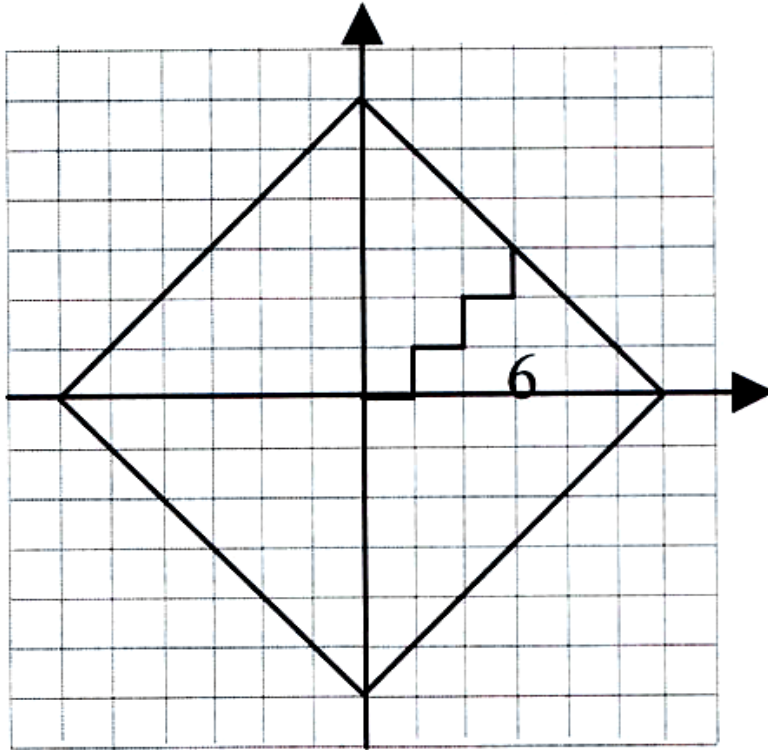
- **Manhattan distance (city block distance):**

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

- A particular case of the Minkowski distance,  $p = 1$
- Measured in “blocks”



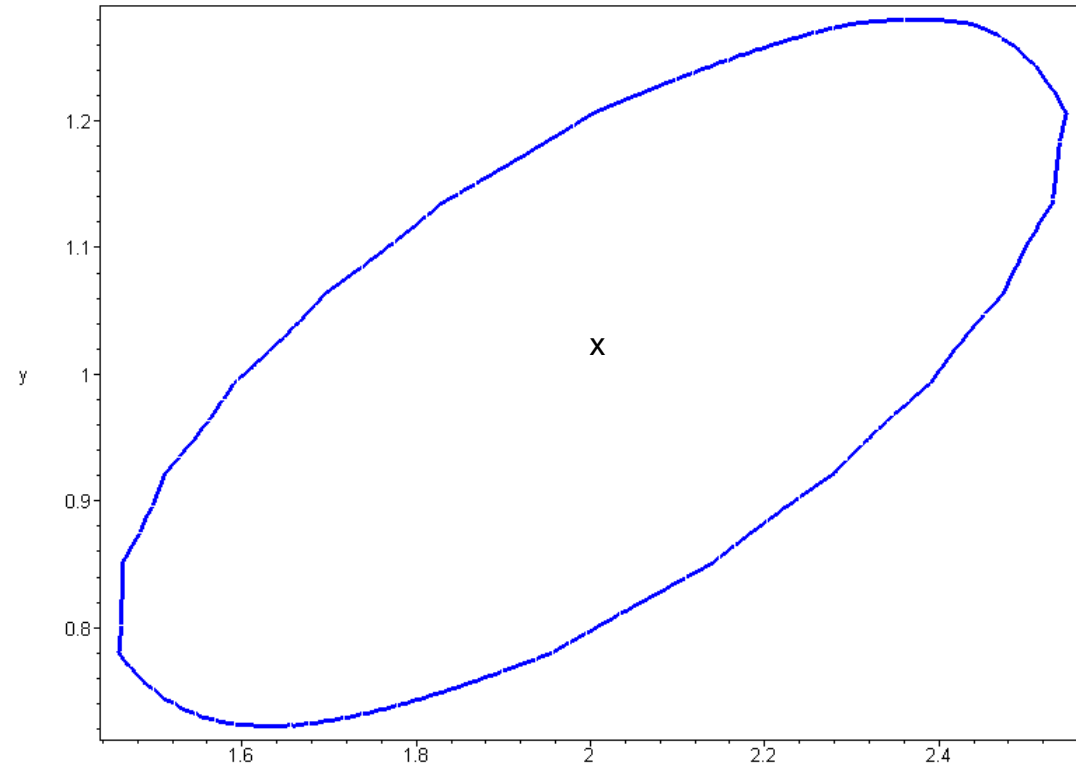
# Comparison: Euclidean vs. Manhattan



# Mahalanobis distance

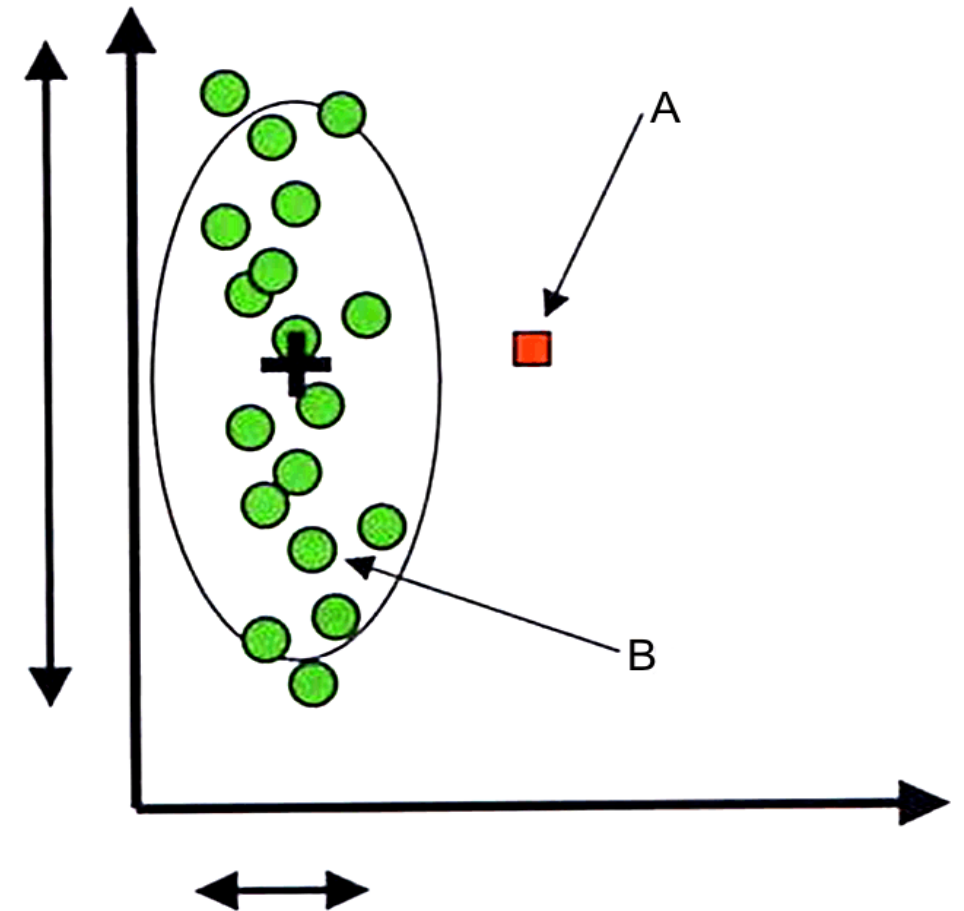
- A positive, definite, symmetric matrix:  $\Sigma$

$$d(\mathbf{x}, \mathbf{y}) = \left[ (\mathbf{x} - \mathbf{y})^t \Sigma^{-1} (\mathbf{x} - \mathbf{y}) \right]^{\frac{1}{2}}$$



$$\Sigma := \begin{bmatrix} 1 & 0.35 \\ 0.35 & 0.2625 \end{bmatrix}$$

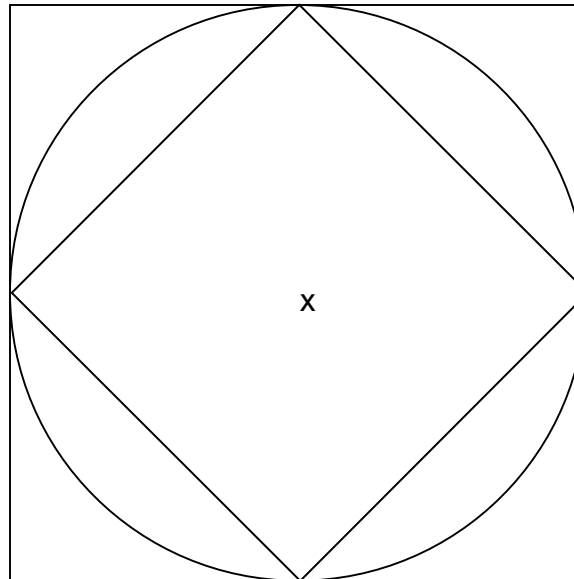
Uncorrelated random variables:





# $L_p$ -norm

- It's the Minkowski distance
- When:
  - $p = 2$ :  $L_2$ -norm or Euclidean distance
  - $p = 1$ :  $L_1$ -norm or Manhattan distance
  - $p = \infty$ :  $L_\infty$ -norm or **sup distance**



- **Chebychev distance:**

- Largest difference between any two components

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

- **Angle distance:**

- Angle between two vectors

$$d(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^t \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

# Correlation distances

*Pearson* correlation distance between  $\mathbf{x}$  and  $\mathbf{y}$

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{\sum_{i=1}^d (x_i - \bar{x})(y_i - \bar{y})}{\left[ \sum_{i=1}^d (x_i - \bar{x})^2 \right]^{1/2} \left[ \sum_{i=1}^d (y_i - \bar{y})^2 \right]^{1/2}}$$

Linear  
correlation

where  $\bar{x} = \sum_{i=1}^d x_i$  and  $\bar{y} = \sum_{i=1}^d y_i$

*Spearman* correlation distance between  $\mathbf{x}$  and  $\mathbf{y}$

- Create two vectors  $\mathbf{a} = [a_1, \dots, a_d]$  and  $\mathbf{b} = [b_1, \dots, b_d]$ ,
- where  $a_i$  and  $b_i$  are the ranks of  $x_i$  and  $y_i$  respectively.

Compute:

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{6 \sum_{i=1}^d (a_i - b_i)^2}{d(d^2 - 1)}$$

Nonlinear  
correlation

# Example

$$\begin{aligned}\mathbf{x} &= [0.1 \quad 0.7 \quad 1.2 \quad 0.5] & \mathbf{y} &= [3.5 \quad 4.2 \quad 5.7 \quad 6.2] \\ \mathbf{a} &= [1 \quad 3 \quad 4 \quad 2] & \mathbf{b} &= [1 \quad 2 \quad 3 \quad 4]\end{aligned}$$

$$\sum_{i=1}^d (\mathbf{a}_i - \mathbf{b}_i)^2 = (\mathbf{a} - \mathbf{b})^t (\mathbf{a} - \mathbf{b}) = 0^2 + 1^2 + 1^2 + 2^2 = 6$$

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{6 \cdot 6}{4(16 - 1)} = 1 - \frac{36}{60} = 1 - 0.6 = 0.4$$

# Jaccard index/distance

Let  $S_1$  and  $S_2$  be two sets, the Jaccard index:

$$J(S_1, S_2) = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$$

where  $0 \leq J(S_1, S_2) \leq 1$

If  $S_1$  and  $S_2$  are both empty:  $J(S_1, S_2) = 1$

Jaccard distance is defined as:

$$d_J = 1 - J(S_1, S_2)$$

Jaccard distance is used for nominal or textual features

For example, phone brand name (Apple, Samsung, Huawei, LG, Nokia, Blackberry), city, province, car make, etc.

# Comparison of Distance Functions

- **Euclidean distance:** Most common distance function used by human beings – it's like an optimal way to “reach something”
- **Angle between vectors:** Considers *only* angle, and *not* magnitude of the vectors  
e.g.  $\mathbf{x} = [1,1]$ ,  $\mathbf{y} = [100,100]$  ( $\cos \pi/2 = 0$ )
- **Orthogonal vectors:** Two vectors are orthogonal iff angle distance is  $\cos 0 = 1$   
e.g.  $\mathbf{x} = [1,1]$ ,  $\mathbf{y} = [-100,100]$
- **Correlation distance:** Looks for similar *variations*, and not for actual numerical values  
e.g.  $\mathbf{x} = [1,2,3,4,5]$ ,  $\mathbf{y} = [10,20,30,40,50]$ ,  $\mathbf{z} = [5,4,3,2,1]$ 
  - $d(\mathbf{x},\mathbf{y})$  will be small, while  $d(\mathbf{x},\mathbf{z})$  will be large
  - Euclidean distance is the opposite

- **Mahalanobis distance:** Takes correlation of r.v. into consideration  
If r.v. are uncorrelated, takes covariances into account
- **Manhattan:** More appropriate for discrete r.v.
- **Minkowski:** A generalization of the Euclidean distance.
- **Chebychev:** Focuses on the *most important* difference.  
e.g.  $\mathbf{x} = [1,2,3,4]$ ,  $\mathbf{y} = [2,3,4,5]$   
 $\mathbf{x} = [1,2,3,4]$ ,  $\mathbf{y} = [1,2,3,6]$

# Edit distance

- When samples are not real-valued vectors, but strings...
- Can use *edit distance*
- In this case, there is no *obvious* measure for the distance
- For example,
- It is not clear whether *abbccc* is closer to *aabbcc* or to *abbcccb*

- **Example:**

*X = excused* is transformed into

*Y = exhausted*

- **Steps:**

- First, substitute *h* for *c*, yielding *X = exhused*
- Second, insert *a*, obtaining *X = exhaused*
- Third, insert *t*, obtaining *X = exhausted*
- Since the cost of each operation is 1, the distance between these two strings is 3



# Edit distance – cont'd

- Given two strings, **x** and **y**.
- The edit distance is based on three operations:
- **Substitution**: A character in **x** is replaced by the corresponding character in **y**.
- **Insertion**: A character in **y** is inserted into **x**.
- **Deletion**: A character in **x** is deleted.
- Another operation:
  - **Interchange**: Exchange two neighbor characters in **x**  
e.g., transform **x** = *abc* into **y** = *cba*
- A dynamic programming algorithm for the edit distance follows...
- It uses a matrix of costs, or “distances”, and
- $\delta(\mathbf{x}[i], \mathbf{y}[j])$ , which is 1 if  $\mathbf{x}[i] = \mathbf{y}[j]$ , and 0 otherwise

# Algorithm Edit Distance

```
C[0,0] ← 0
for  $i \leftarrow 0$  to  $m$                                 // length of x
    C[ $i$ ,0] ←  $i$ ;
endfor
for  $j \leftarrow 0$  to  $n$                                 // length of y
    C[0, $j$ ] ←  $j$ ;
endfor
for  $i \leftarrow 0$  to  $m$ 
    for  $j \leftarrow 0$  to  $n$ 
        C[ $i$ , $j$ ] ← min{ C[ $i-1$ , $j$ ]+1, C[ $i$ ,  $j-1$ ]+1,
                        C[ $i-1$ , $j-1$ ]+1-  $\delta(\mathbf{x}[i],\mathbf{y}[j])$  }
        Insertion ← C[ $i-1$ , $j$ ]+1
        Deletion ← C[ $i$ ,  $j-1$ ]+1
        No change / exchange ← C[ $i-1$ , $j-1$ ]+1-  $\delta(\mathbf{x}[i],\mathbf{y}[j])$ 
    endfor
endfor
return C[ $m,n$ ]
```

$$\delta(\mathbf{x}[i],\mathbf{y}[j]) = \begin{cases} 1 & \text{if } \mathbf{x}[i] = \mathbf{y}[j] \\ 0 & \text{otherwise} \end{cases}$$

Complexity:  $O(mn)$

# Example

C	e x h a u s t e d								
	1	2	3	4	5	6	7	8	9
e	1	0	1	2	3	4	5	6	7
x	2	1	0	1	2	3	4	5	6
c	3	2	1	1	2	3	4	5	6
u	4	3	2	2	2	2	3	4	5
s	5	4	3	3	3	3	2	3	4
e	6	5	4	4	4	4	3	3	4
d	7	6	5	5	5	5	4	4	3

- Algorithm EditDistance runs in  $O(nm)$
- Function  $\delta$  can be generalized to:
  - Other than 0-1 cost function
  - Can consider different costs for different symbols
- Generalization includes protein/DNA sequence alignment

**Algorithm** *EditDistance*( $X, Y$ )

**Input:** Strings  $X$  and  $Y$  of length  $n$  and  $m$  respectively

**Output:** Array  $C$  containing prefix edit distances

$C[0,0] \leftarrow 0$

**for**  $i \leftarrow 0$  **to**  $m$  // length of  $X$

$C[i,0] \leftarrow i$

**for**  $j = 0$  **to**  $n$  // length of  $Y$

$C[0,j] \leftarrow j$

**for**  $i = 0$  **to**  $m$

**for**  $j = 0$  **to**  $n$

$C[i,j] = \min\{C[i-1,j]+1, C[i, j-1]+1, \\ C[i-1,j-1]+1- \delta(X_i, Y_j) \}$

**return**  $C[m,n]$

# Bio-sequence alignment

- General case: global alignment
- Bio-sequences can represent:
  - DNA, RNA, proteins
- Uses a scoring system that assigns:
  - score values for **each** pair of symbols
  - e.g. nucleic acid or amino acid pairs
  - penalty values to single gaps
- Score = sum of pair scores – sum of gap penalties

Example:  $X = \text{TGATAGCCAG}$      $Y = \text{AGAGCA}$

T - G A T A G C C A G

- A G A G - - C - A -

Score =  $-4 - 4 + 6 + 6 + 2 - 4 - 4 + 6 - 4 + 6 - 4 = 2$

Score for optimal alignment of  $X_i$  and  $Y_j$  given by:

$$c_{ij} = \max \left\{ c_{i-1,j-1} + s(x_i, y_j), \right. \\ \left. \max_{k \geq 1} (c_{i-k,j} - w_k), \right. \\ \left. \max_{t \geq 1} (c_{i,j-t} - w_t) \right\} \quad (1)$$

Score for optimal alignment of  $X$  and  $Y = c_{mn}$

**Algorithm** SeqAlign( $X, Y$ )

$c_{i0} \leftarrow -w_i, \quad 1 \leq i \leq m$

$c_{0j} \leftarrow -w_j, \quad 1 \leq j \leq n$

**for**  $i = 1$  **to**  $m$

**for**  $j = 1$  **to**  $n$

        Compute  $c_{ij}$  as in (1)

        Update “traceback” matrix  $D$  (with  $\leftarrow \uparrow \nwarrow$ )

**return**  $c_{mn}$

Worst-case running time:  $O(nm)$

	A	C	G	T
A	6	2	2	2
C	2	6	2	2
G	2	2	6	2
T	2	2	2	6

# DNA sequence alignment – example

	$\lambda$	T	G	A	T	A	G	C	C	A	G
$\lambda$	0	-4	-8	-12	-16	-20	-24	-28	-32	-36	-40
A	-4	2	-2	-2	-6	-10	-14	-18	-22	-26	-30
G	-8	-2	8	4	0	-4	-4	-8	-12	-16	-20
A	-12	-6	4	14	10	6	2	-2	-6	-6	-10
G	-16	-10	0	10	16	12	12	8	4	0	0
C	-20	-14	-4	2	12	18	14	18	14	10	6
A	-24	-18	-8	2	8	18	20	16	20	20	16

Optimal alignment ( $\leftarrow$ ):

A G A - - G C - A -  
T G A T A G C C A G

$$C_{mn} = 2+6+6-4-4+6+6-4+6-4 = 16$$

	A	C	G	T
A	6	2	2	2
C	2	6	2	2
G	2	2	6	2
T	2	2	2	6

- Dynamic programming can obtain more than one solution
- e.g., following the pink path
- Algorithm can be extended to protein sequence alignment
  - Considers 20 amino acids and other score matrices