

Data Mining

Naïve Bayes/KNN

<https://data-mining.github.io/winter-2026/>

CS 453/553 – Winter 2026

Yu Wang, Ph.D.

Assistant Professor

Computer Science

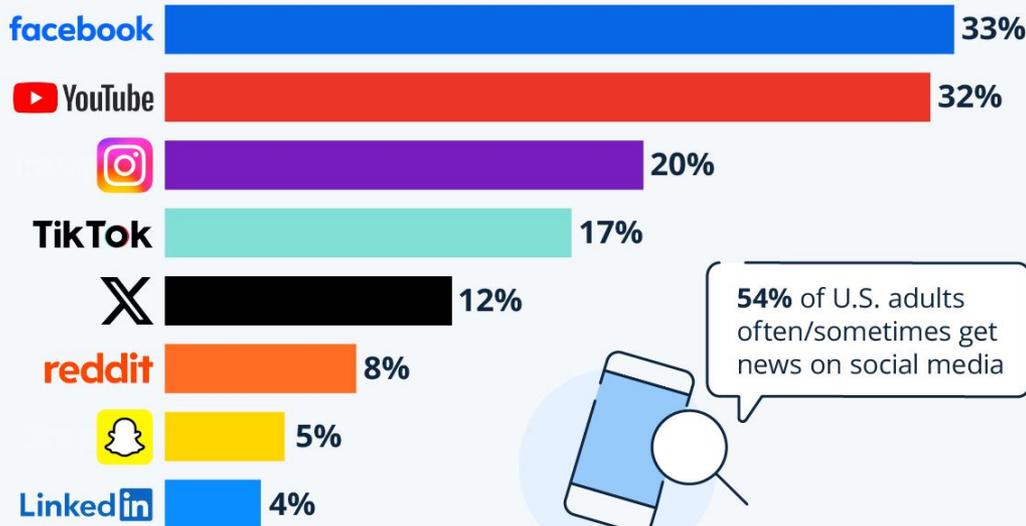
University of Oregon



Motivation – What is Classification and Why?

54% of Americans Get (Mis)informed on Social Media

Share of U.S. adults who regularly get news on the following social media platforms



54% of U.S. adults often/sometimes get news on social media

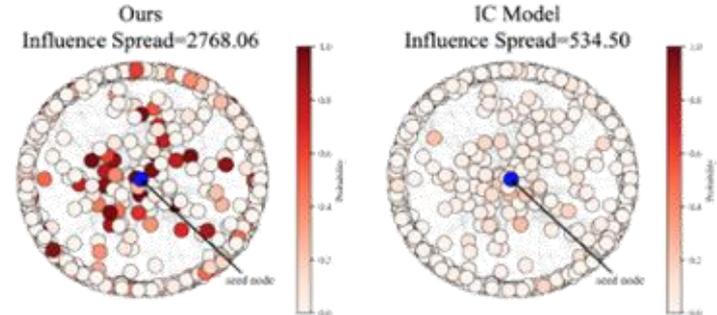
10,658 U.S. adults surveyed Jul-Aug 2024
Source: Pew Research Center



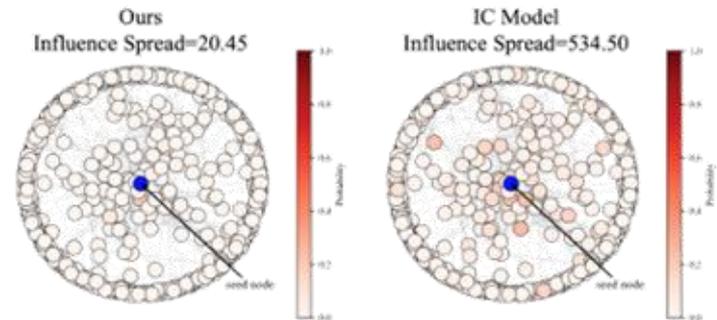
statista

Misinformation Detection

Text: "Breaking: NASA confirms first-ever human colony on Mars will begin next year — tickets for civilians already being sold out in minutes!"

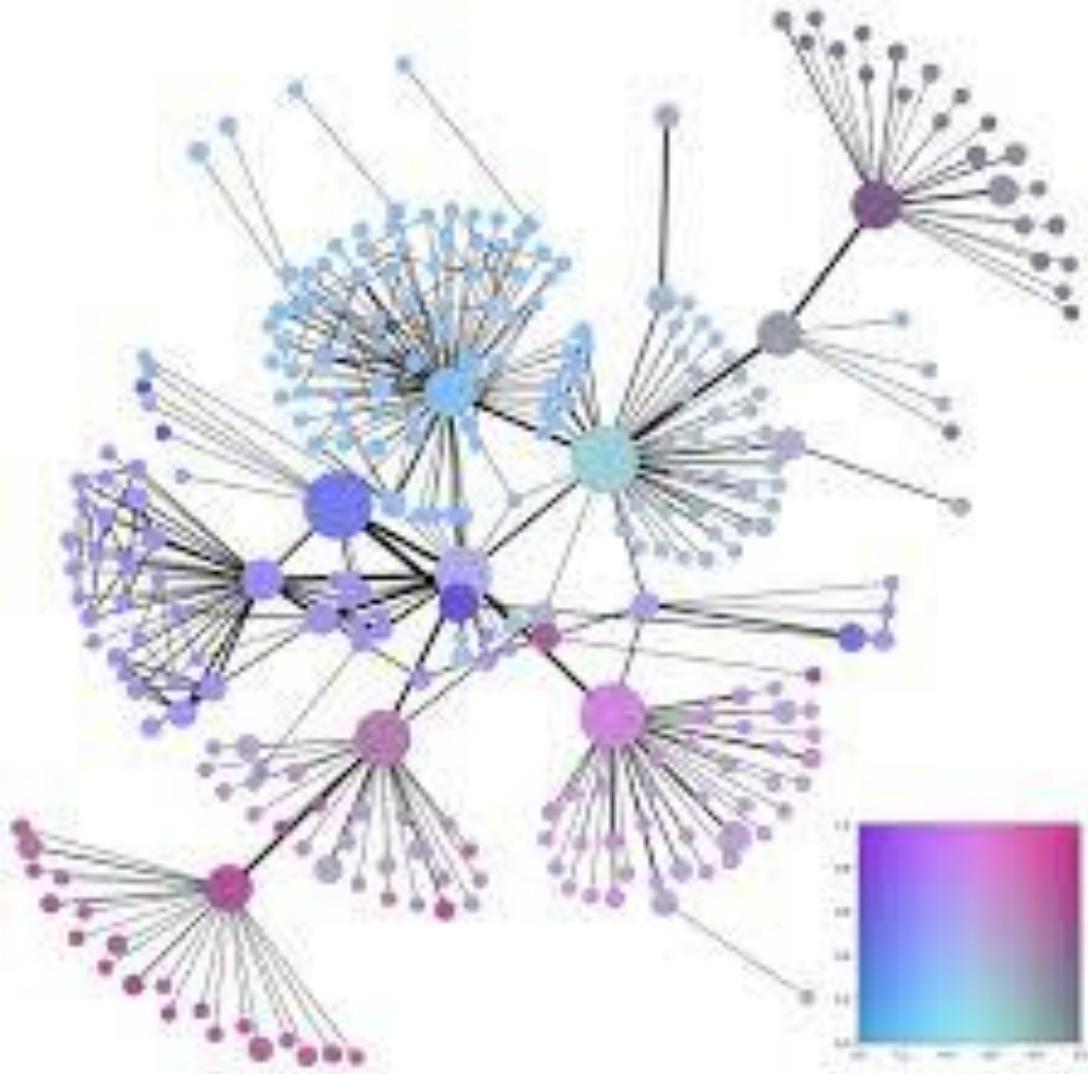


Text: "Today I bought a new pencil."

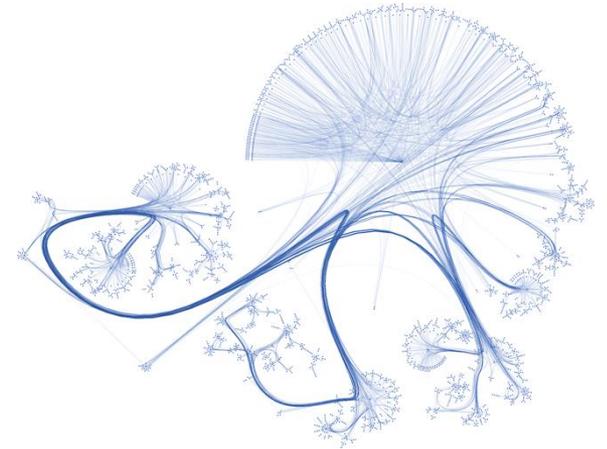




Motivation – What is Classification and Why?



Paper Classification





Some Basic Methods

- **Naïve Bayes Classifier**
- **Nearest Neighbor Classifiers**
- **Decision-Tree**
- **Advanced Methods:**
 - **Neural Network**
 - **Geometric Neural Networks**





Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

Whether Evade?

Yes?

No?

$P(\text{Yes} \mid X)$

$P(\text{No} \mid X)$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$$P(\text{No}) = 7/10,$$
$$P(\text{Yes}) = 3/10$$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$P(X|\text{Yes})$

$P(X|\text{No})$

$P(X_{re}, X_d, X_{in} | \text{Yes})$

$P(X_{re}, X_d, X_{in} | \text{No})$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$$P(X_{re}, X_d, X_{in} | \text{Yes}) = P(X_{re} | \text{Yes}) \\ P(X_d | \text{Yes})P(X_{in} | \text{Yes})$$

$$P(X_{re}, X_d, X_{in} | \text{No}) = P(X_{re} | \text{No}) \\ P(X_d | \text{No})P(X_{in} | \text{No})$$

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$P(X | \text{Yes}) =$

$P(\text{Refund} = \text{No} | \text{Yes}) \times$
 $P(\text{Divorced} | \text{Yes}) \times$
 $P(\text{Income} = 120\text{K} | \text{Yes})$

$P(X | \text{No}) =$

$P(\text{Refund} = \text{No} | \text{No}) \times$
 $P(\text{Divorced} | \text{No}) \times$
 $P(\text{Income} = 120\text{K} | \text{No})$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

For categorical attributes:

$$P(X_i = c | y) = n_c / n$$

- where $|X_i = c|$ is number of instances having attribute value $X_i = c$ and belonging to class y
- Examples:

$$P(\text{Status}=\text{Married}|\text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes}|\text{Yes})=0$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$P(X | \text{Yes}) =$

$P(\text{Refund} = \text{No} | \text{Yes}) \times$

$P(\text{Divorced} | \text{Yes}) \times$

$P(\text{Income} = 120\text{K} | \text{Yes})$

$P(X | \text{No}) =$

$P(\text{Refund} = \text{No} | \text{No}) \times$

$P(\text{Divorced} | \text{No}) \times$

$P(\text{Income} = 120\text{K} | \text{No})$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$

$$P(\text{Yes} | X) = \frac{P(X | \text{Yes})P(\text{Yes})}{P(X)}$$

$$P(\text{No} | X) = \frac{P(X | \text{No})P(\text{No})}{P(X)}$$

$$P(X|Y = \text{No}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X-\mu)^2}{2\sigma^2}}$$

If Class=No

- ◆ sample mean = 110
- ◆ sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi(54.54)}} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



Naïve Bayes Classifier

Given a Test Record:

$$X = (\text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120\text{K})$$

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} \mid \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} \mid \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} \mid \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

- $P(X \mid \text{No}) = P(\text{Refund}=\text{No} \mid \text{No})$
 $\times P(\text{Divorced} \mid \text{No})$
 $\times P(\text{Income}=120\text{K} \mid \text{No})$
 $= 4/7 \times 1/7 \times 0.0072 = 0.0006$
- $P(X \mid \text{Yes}) = P(\text{Refund}=\text{No} \mid \text{Yes})$
 $\times P(\text{Divorced} \mid \text{Yes})$
 $\times P(\text{Income}=120\text{K} \mid \text{Yes})$
 $= 1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times 10^{-10}$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$

\Rightarrow Class = No



Naïve Bayes Classifier

Even in absence of information about any attributes, we can use Apriori Probabilities of Class Variable:

Naïve Bayes Classifier:

$$P(\text{Refund} = \text{Yes} \mid \text{No}) = 3/7$$

$$P(\text{Refund} = \text{No} \mid \text{No}) = 4/7$$

$$P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0$$

$$P(\text{Refund} = \text{No} \mid \text{Yes}) = 1$$

$$P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7$$

$$P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7$$

$$P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3$$

$$P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3$$

$$P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0$$

For Taxable Income:

If class = No: sample mean = 110

sample variance = 2975

If class = Yes: sample mean = 90

sample variance = 25

$$P(\text{Yes}) = 3/10$$

$$P(\text{No}) = 7/10$$

If we only know that marital status is Divorced, then:

$$P(\text{Yes} \mid \text{Divorced}) = 1/3 \times 3/10 / P(\text{Divorced})$$

$$P(\text{No} \mid \text{Divorced}) = 1/7 \times 7/10 / P(\text{Divorced})$$

If we also know that Refund = No, then

$$P(\text{Yes} \mid \text{Refund} = \text{No}, \text{Divorced}) = 1 \times 1/3 \times 3/10 / P(\text{Divorced}, \text{Refund} = \text{No})$$

$$P(\text{No} \mid \text{Refund} = \text{No}, \text{Divorced}) = 4/7 \times 1/7 \times 7/10 / P(\text{Divorced}, \text{Refund} = \text{No})$$

If we also know that Taxable Income = 120, then

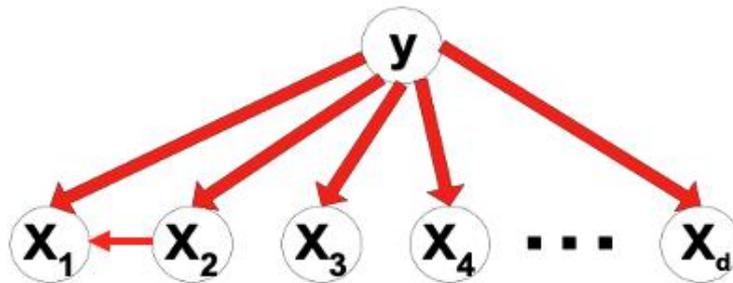
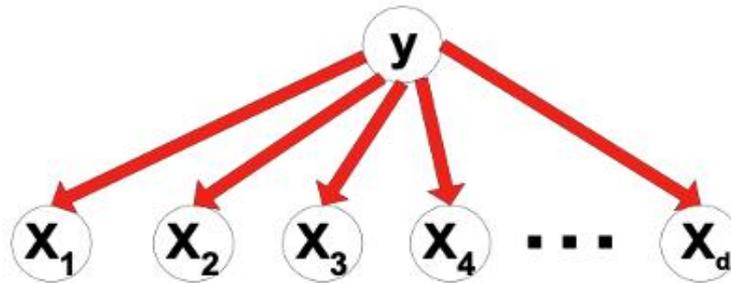
$$P(\text{Yes} \mid \text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120) = 1.2 \times 10^{-9} \times 1 \times 1/3 \times 3/10 / P(\text{Divorced}, \text{Refund} = \text{No}, \text{Income} = 120)$$

$$P(\text{No} \mid \text{Refund} = \text{No}, \text{Divorced}, \text{Income} = 120) = 0.0072 \times 4/7 \times 1/7 \times 7/10 / P(\text{Divorced}, \text{Refund} = \text{No}, \text{Income} = 120)$$



Naïve Bayes Classifier

- Assume independence among attributes X_i when class is given:
 - $P(X_1, X_2, \dots, X_d | Y_j) = P(X_1 | Y_j) P(X_2 | Y_j) \dots P(X_d | Y_j)$





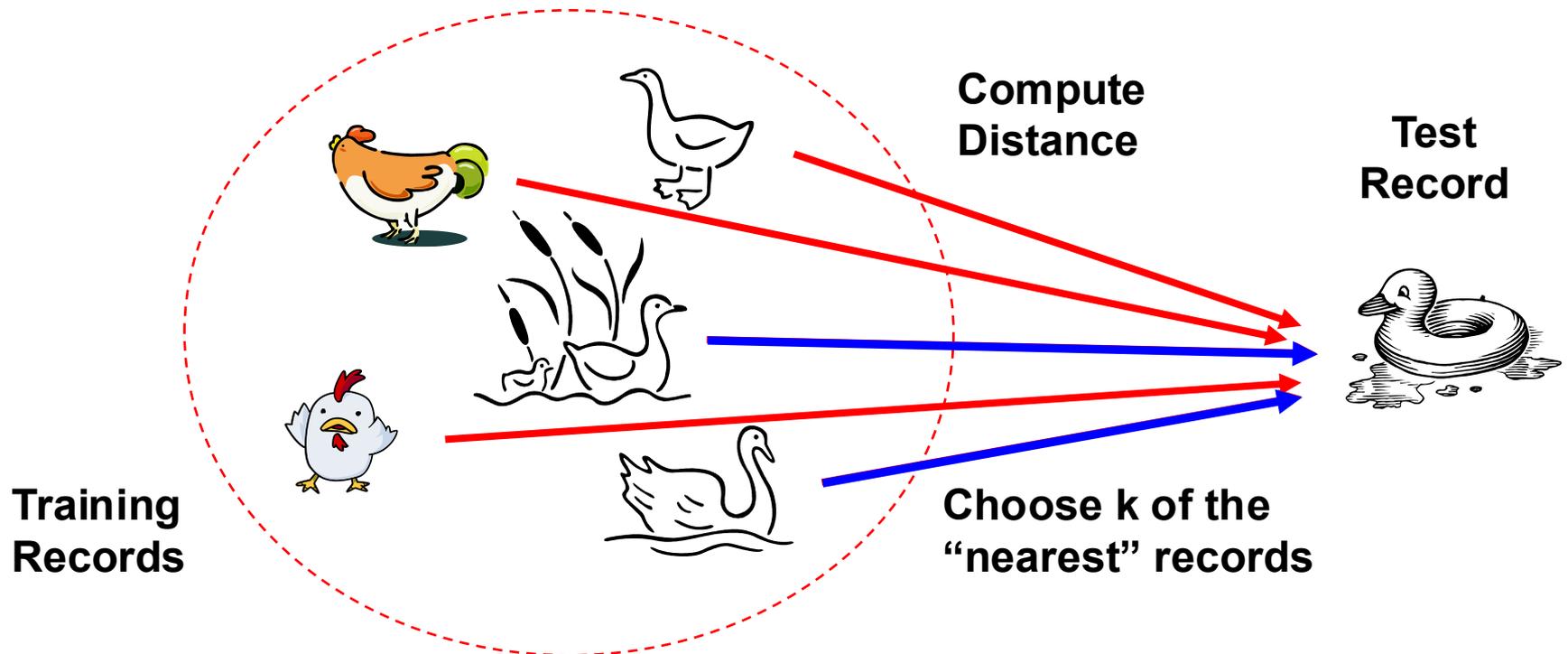
Some Basic Methods

- **Naïve Bayes Classifier**
- **Nearest Neighbor Classifiers**
- **Decision-Tree**
- **Advanced Methods:**
 - **Neural Network**
 - **Geometric Neural Networks**



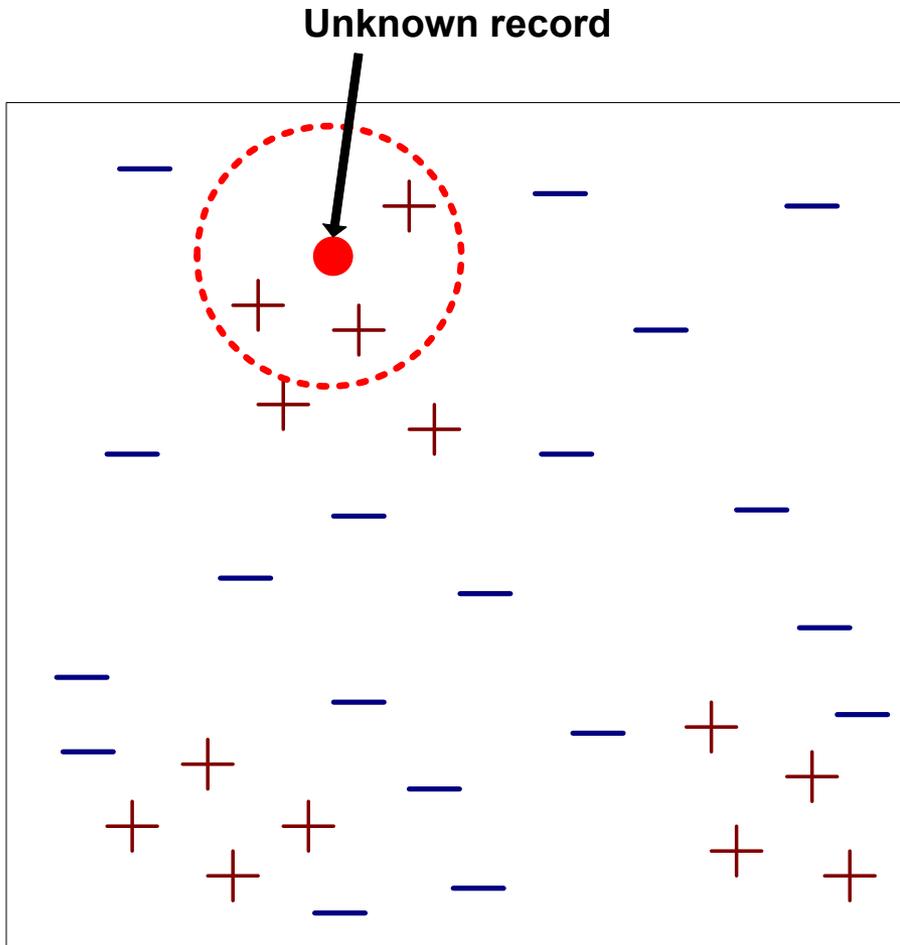
Nearest Neighbor Classifiers

- Basic idea:
 - If it walks like a duck,
quacks like a duck, then it's probably a duck





Nearest Neighbor Classifiers



What is training and testing data?

- Requires the following:
 - A set of labeled records
 - Proximity (Distance/Similarity) metric between a pair of records
 - e.g., Euclidean distance
 - The value of k , the number of nearest neighbors to retrieve
 - class labels of K nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

Or you can do $w = 1/d^2$



Nearest Neighbor Classifiers – Real Estate Example

Commercial

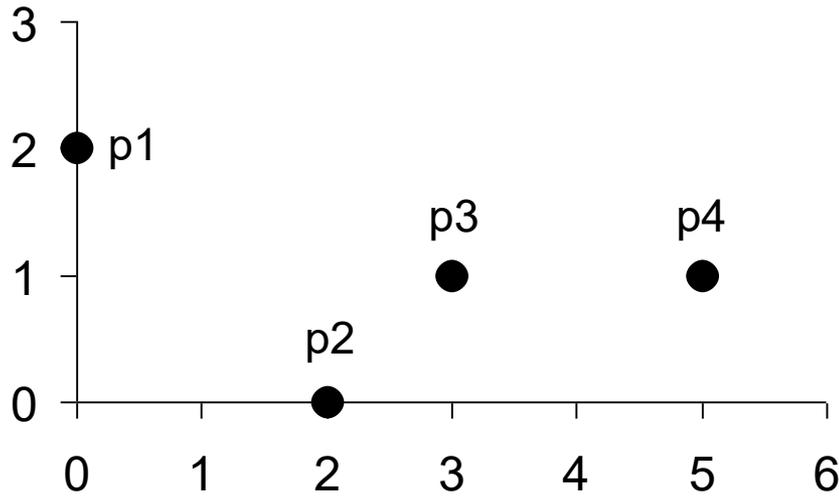


Residential





Nearest Neighbor Classifiers – Distance/Similarity Metric



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0



Nearest Neighbor Classifiers – Paper Classification



KW1 KW2 KW4
[0, 1, 1, 1]



[0, 1, 1, 0]



[1, 0, 0, 0]





Nearest Neighbor Classifiers – Paper Classification



KW1 KW2 KW4
[0, 1, 1, 1]

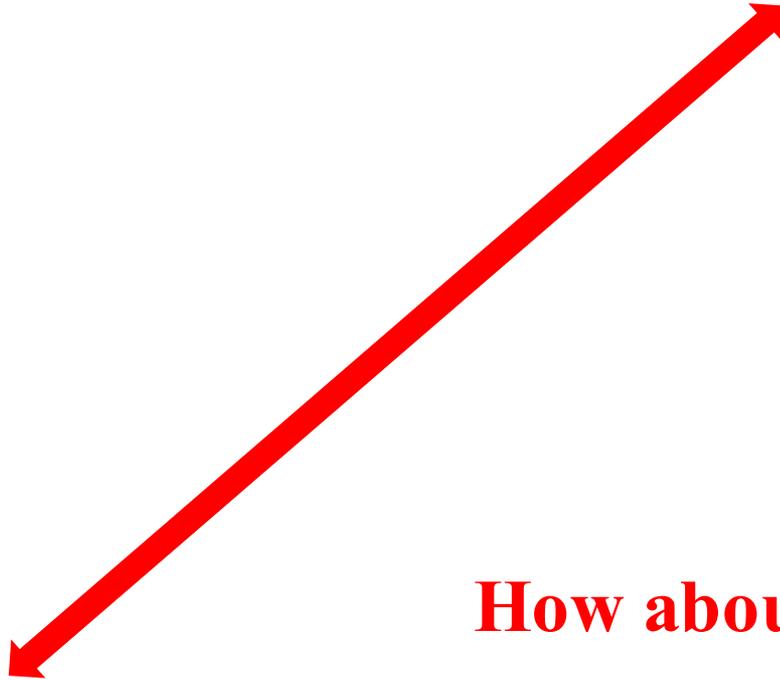
[0, 0, 0, 1]



[0, 1, 1, 0]



[1, 0, 0, 0]



How about these two?



Nearest Neighbor Classifiers – Distance Metric

You need to think about the distance metric!

Inner Product	Minkowski	Intersection	Entropy	χ^2 Family	Fidelity (Squared-Chord)	String Rearrangement	String Similarity
Subsec. 2.1	Subsec. 2.2	Subsec. 2.3	Subsec. 2.4	Subsec. 2.5	Subsec. 2.6	Subsec. 2.7	Subsec. 2.7
Inner Product Cosine Angular Jaccard Dice	Euclidean (Euclidean) ² L_1 L_p L_∞	Intersection Wave Hedges Sørensen Kulczynski Jaccard	Kullback-Leibler J-Divergence K-Divergence Topsøe Jensen-Shannon Jensen Diff. SED	Pearson Neyman Add. Sym. χ^2 Spearman Squared χ^2 Divergence Clark Mahalanobis	Fidelity Bhattacharyya Hellinger Matusia Squared-Chord	Hamming Levenshtein Swap Interchange Parallel-Interchange	LCS Jaro String N-Grams

Table 1. Similarity or distance measures appearing in this survey, categorized by measures families.

<https://arxiv.org/pdf/2408.07706>



Nearest Neighbor Classifiers – Distance Metric

**Cosine
Similarity**

$$\text{sim}_{\text{Cos}}(P, Q) = \frac{\langle P, Q \rangle}{\|P\| \|Q\|}$$

**Jaccard
Similarity**

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

<https://arxiv.org/pdf/2408.07706>



Nearest Neighbor Classifiers – Paper Classification



KW1 KW2 KW4
[0, 1, 1, 1]

[0, 0, 0, 1]



$$\frac{0 * 1 + 1 * 1 + 1 * 1 + 1 * 0}{\sqrt{0 + 1 + 1 + 1} \sqrt{0 + 1 + 1 + 0}} = \frac{2}{\sqrt{6}}$$



[0, 1, 1, 0]

0

**Cosine
Similarity**

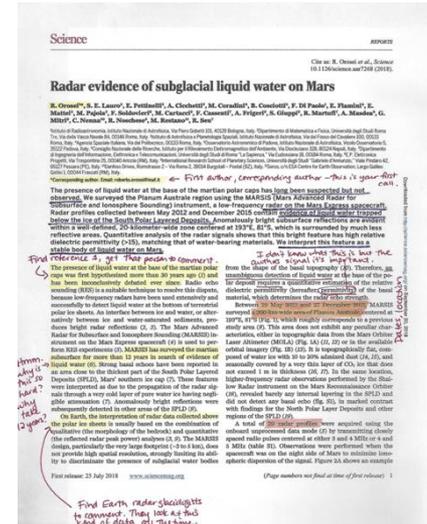
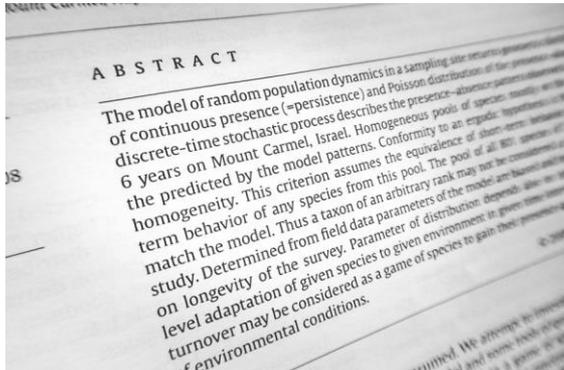


[1, 0, 0, 0]



Nearest Neighbor Classifiers – Feature Preprocessing

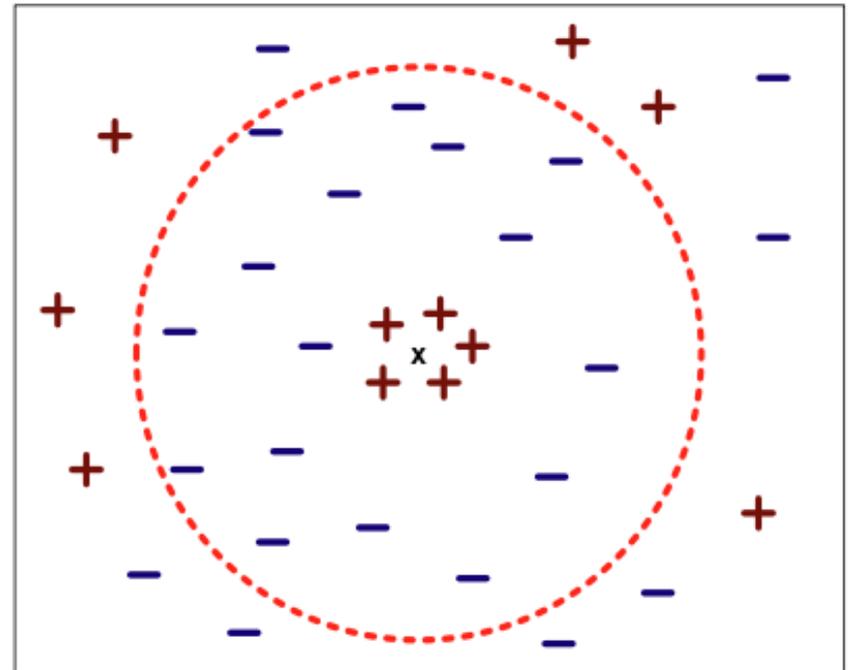
- Data Preprocessing is often required
 - Different attributes have different scales
 - height of a person may vary from 1.5m to 1.8m
 - weight of a person may vary from 90lb to 300lb
 - income of a person may vary from \$10K to \$1M
 - Different instances may have different background





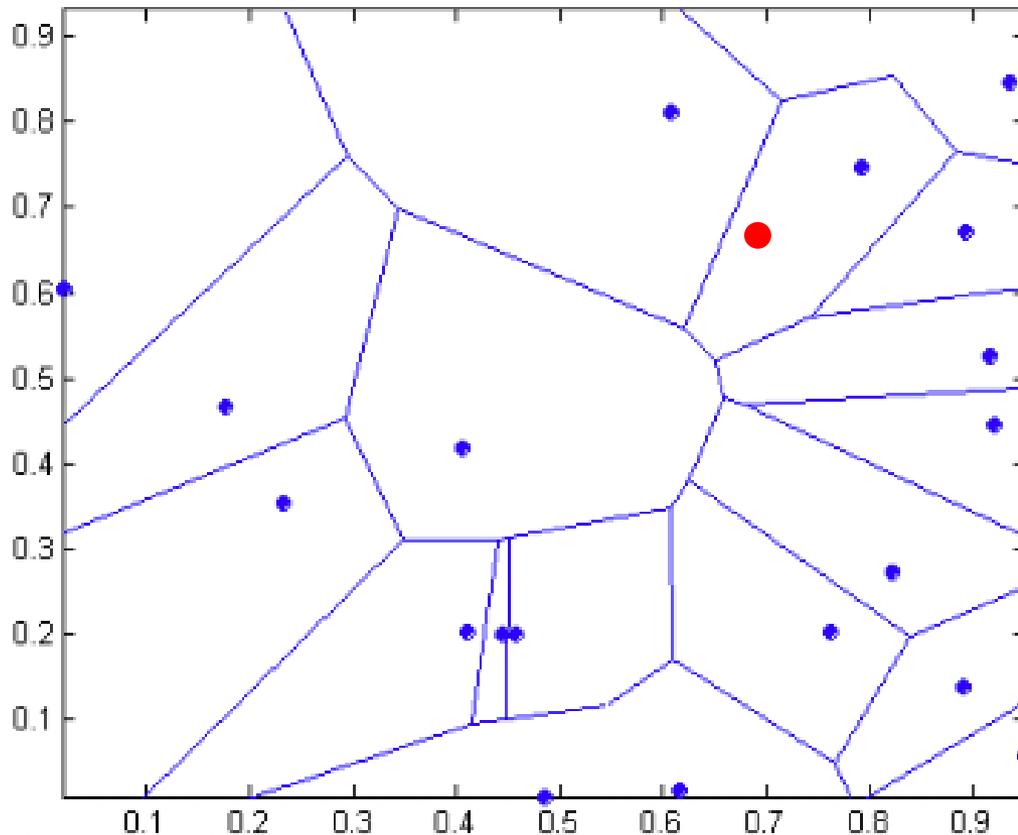
Nearest Neighbor Classifiers – Hyperparameter

- **Choosing the value of k**
- **K too small?**
- **K too large?**





1-nn decision boundary is a Voronoi Diagram



Assuming N Training points

D Dimension

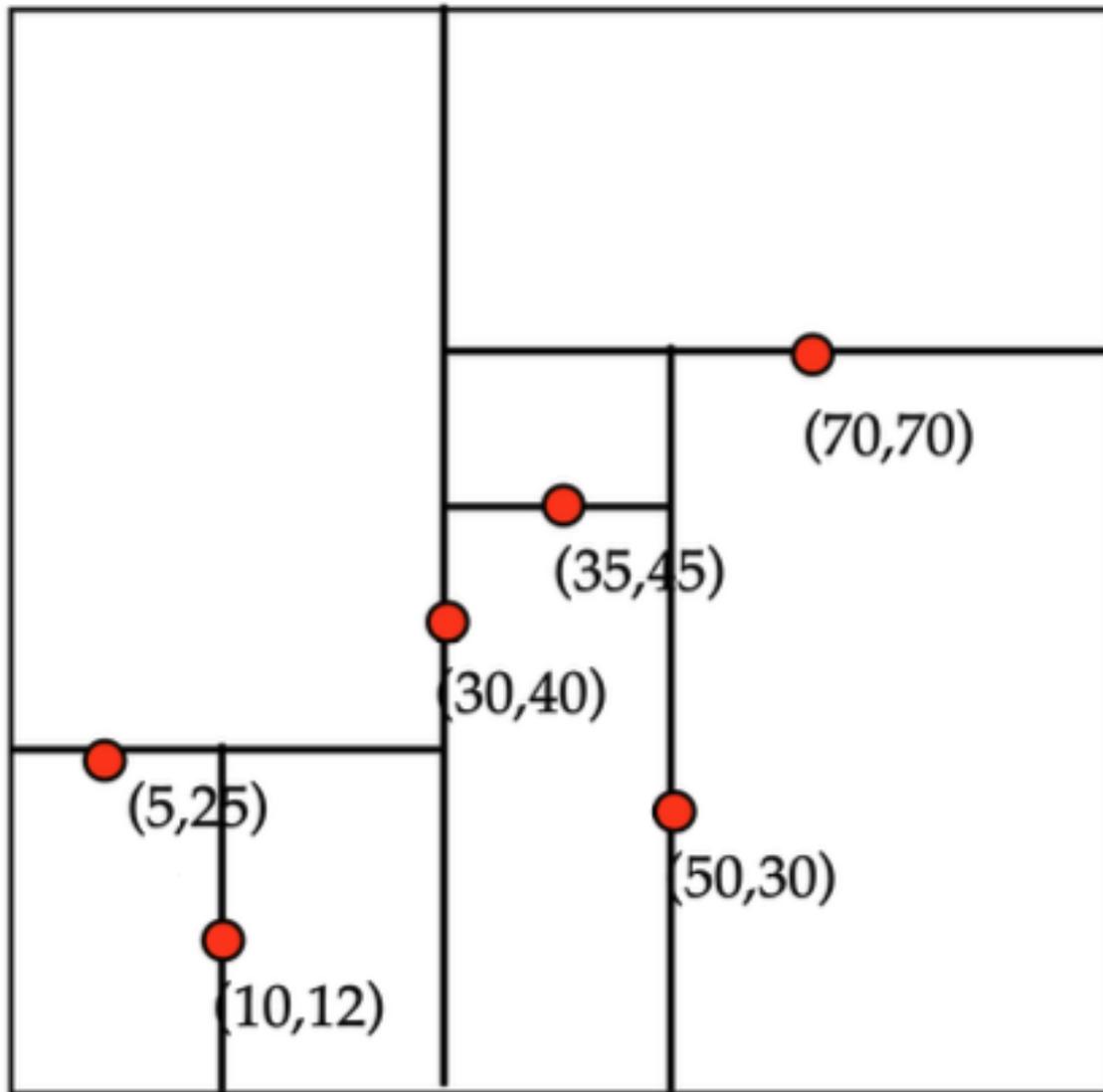
What is the **time complexity** of classifying that red point?

$O(ND)$

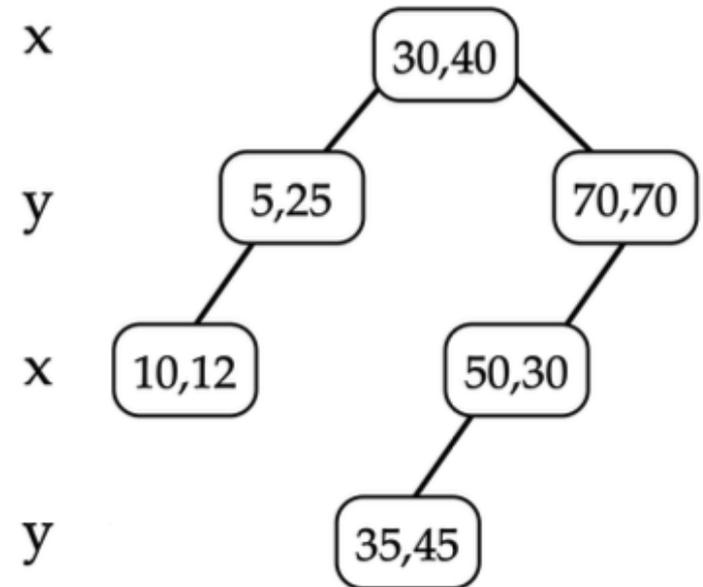
Everytime, we need to compare with every training points



Nearest Neighbor Classifiers – Efficiency



insert: (30,40), (5,25), (10,12), (70,70), (50,30), (35,45)

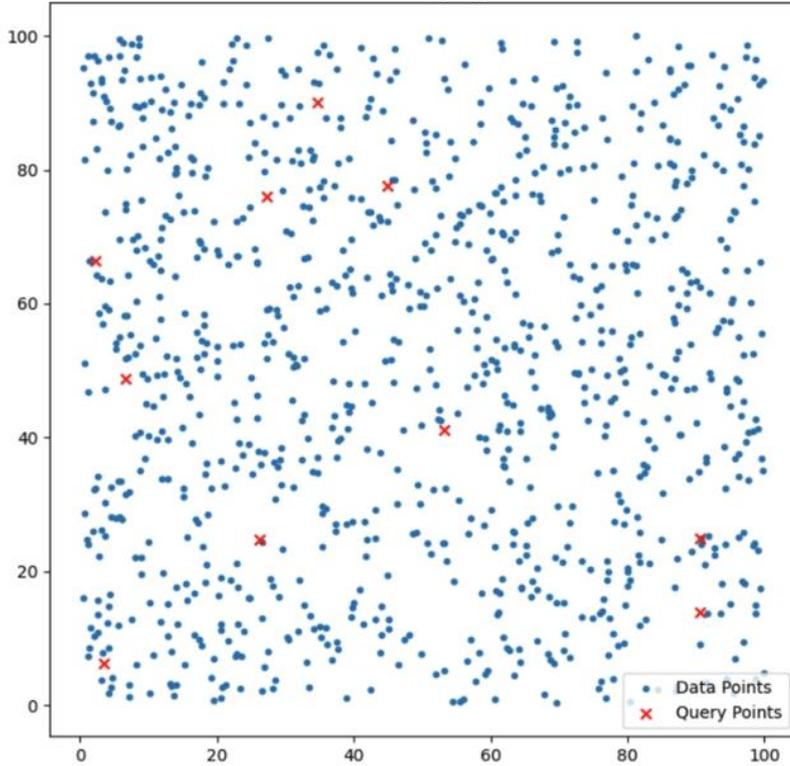


What is the **time complexity** of classifying that red point?

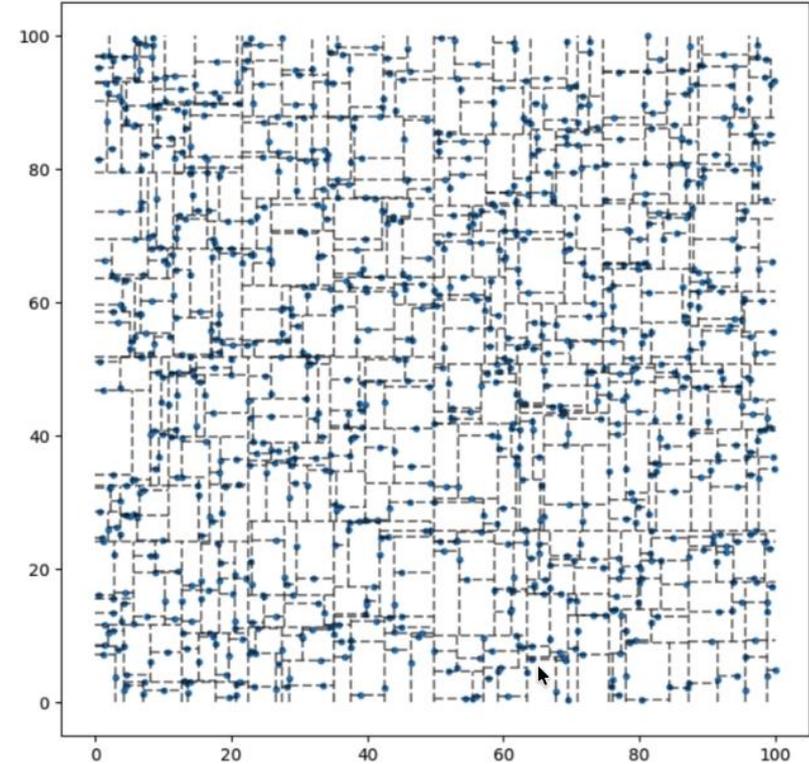


Nearest Neighbor Classifiers – Efficiency

Dataset with Query Points



KD-Tree Partitioning



```
# Generate synthetic 2D dataset
np.random.seed(42)
n_samples = 1000 # Number of data points
data = np.random.rand(n_samples, 2) * 100 # 2D points in a 100x100 area

# Generate some random query points
n_queries = 10
queries = np.random.rand(n_queries, 2) * 100
```

```
Brute-force avg time: 0.175 ms
KD-Tree avg time: 0.059 ms
```



Nearest Neighbor Classifiers – Efficiency

Meta Our approach ▾ Research ▾ Product experiences ▾ Llama Blog

TOOLS

Faiss

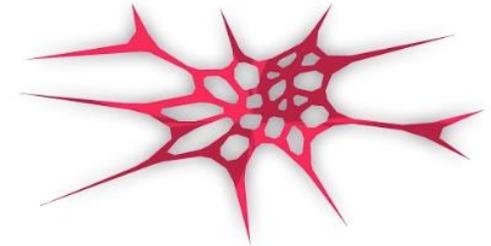
Faiss (Facebook AI Similarity Search) is a library that allows developers to quickly search for embeddings of multimedia documents that are similar to each other. It solves limitations of traditional query search engines that are optimized for hash-based searches, and provides more scalable similarity search functions.

Efficient similarity search

With Faiss, developers can search multimedia documents in ways that are inefficient or impossible with standard database engines (SQL). It includes nearest-neighbor search implementations for million-to-billion-scale datasets that optimize the memory-speed-accuracy tradeoff. Faiss aims to offer state-of-the-art performance for all operating points.

Faiss contains algorithms that search in sets of vectors of any size, and also contains supporting code for evaluation and parameter tuning. Some of its most useful algorithms are implemented on the GPU. Faiss is implemented in C++, with an optional Python interface and GPU support via CUDA.

FAISS
Scalable Search With Facebook AI



```
# Generate synthetic 2D dataset
np.random.seed(42)
n_samples = 100000
data = np.random.rand(n_samples, 20) * 100

n_queries = 10
queries = np.random.rand(n_queries, 20) * 100
```

```
Brute-force avg time: 8.112 ms
KD-Tree avg time:    4.666 ms
FAISS (FlatL2) avg:  0.253 ms
```