

## Predict Party Affiliation of Congressmen

### GROUP MEMBERS:

SPANDAN MAITY, CEMK, REG NO- 151070110050 OF 2015-2016

DEBANJAN DAS, CEMK, REG NO- 151070110021 OF 2015-2016

*y. Pradyumn Choudhury*



# CONTENTS

1. ACKNOWLEDGEMENT
2. DESCRIPTION OF THE DATA
3. GRAPHICAL ANALYSIS
4. DATA ANALYSIS
5. CODE
6. RESULT
7. FINAL COLUMN EXCLUSION
8. EXPERIMENT
9. FUTURE IMPROVEMENT
10. CERTIFICATE

*y Pro Chowdhury*



## ACKNOWLEDGEMENT

I take this opportunity to express my profound gratitude and deep regards to my faculty **Mr. Titas Roy Chowdhury** for his exemplary guidance, monitoring and constant encouragement throughout the course of this project. The blessing, help and guidance given by him/her time to time shall carry me a long way in the journey of life on which I am about to embark.

I am obliged to my project team members for the valuable information provided by them in their respective fields. I am grateful for their cooperation during the period of my assignment.



## DESCRIPTION OF THE DATA

- The voting record of the current House of Representatives on 4 important votes in 2003.
- **Content :**
- State: Name of the States of the corresponding Congressmen.
- Name: Name of the Congressmen.
- Party: Name of party of the corresponding congressmen.
- Vote1: Vote given in 1<sup>st</sup> stage ("Y" for yes, "N" for vote against the topic, "-" for no vote)
- Vote2: Vote given in 2<sup>nd</sup> stage ("Y" for yes, "N" for vote against the topic, "-" for no vote)
- Vote3: Vote given in 3<sup>rd</sup> stage ("Y" for yes, "N" for vote against the topic, "-" for no vote)
- Vote4: Vote given in 4<sup>th</sup> stage ("Y" for yes, "N" for vote against the topic, "-" for no vote)
- **Goal :** To use Naive Bayes and nearest neighbors learning methods to predict the party affiliation of congressmen from their voting records.

*y. Pradyumn Chowdhury*



- **COLUMN NAME DESCRIPTION:**

<b><u>Columns</u></b>	<b><u>Types</u></b>	<b><u>Variable</u></b>
State	object	Continuous
Name	object	Continuous
Party	object	Categorical(D , R)
Vote1	object	Categorical(Y , N , - )
Vote2	object	Categorical(Y , N , - )
Vote3	object	Categorical(Y , N , - )
Vote4	object	Categorical(Y , N , - )

- There are 434 rows and 7 columns in the Given Data.

*y. Prashant Chowdhury*



## GRAPHICAL ANALYSIS:

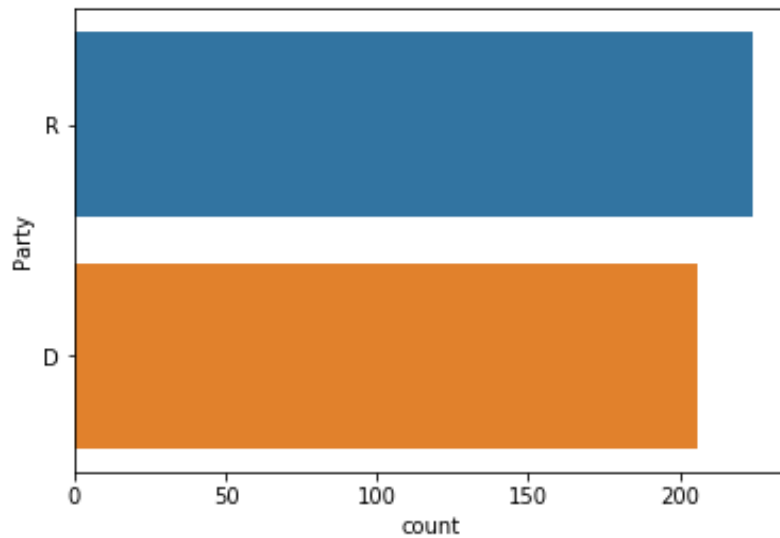
*y/Prof Chowdhury*



## PARTY:

### Code:

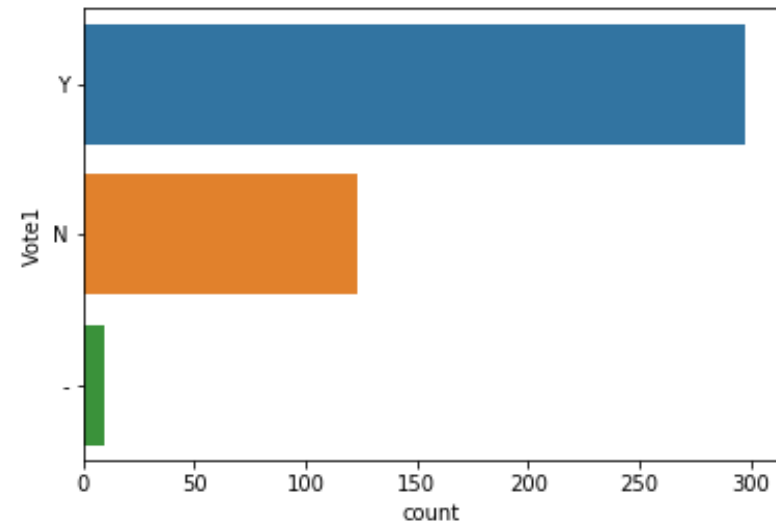
```
sns.countplot(y="Party",data=vote)
```



## VOTE1:

### Code:

```
sns.countplot(y="Vote1",data=vote)
```



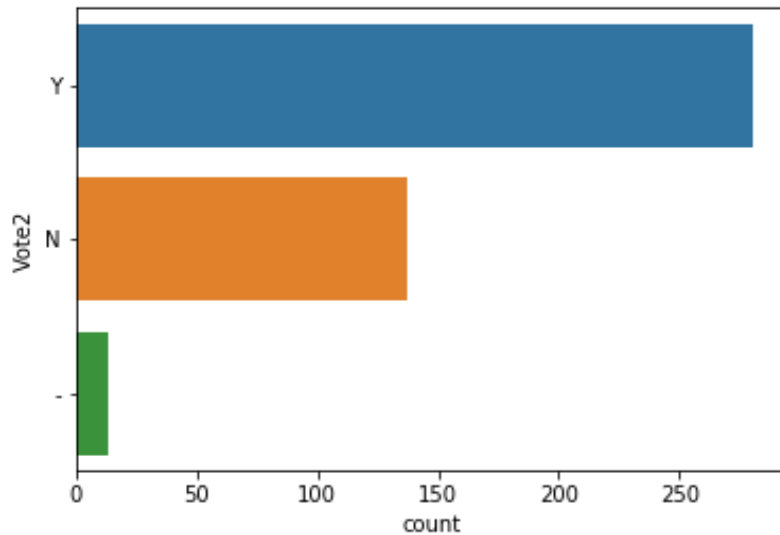
*y. Prasad Chowdhury*



## VOTE2:

### Code:

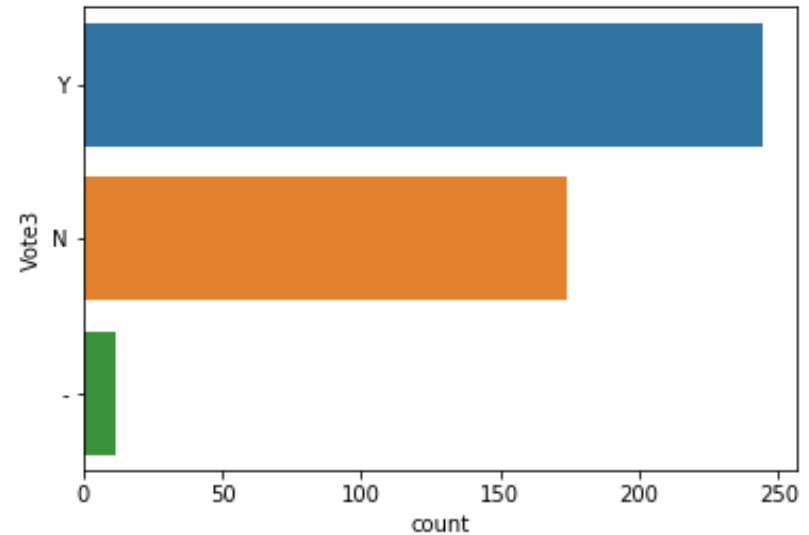
```
sns.countplot(y="Vote2",data=vote)
```



## VOTE3:

### Code:

```
sns.countplot(y="Vote3",data=vote)
```



*y. Prasad Chowdhury*

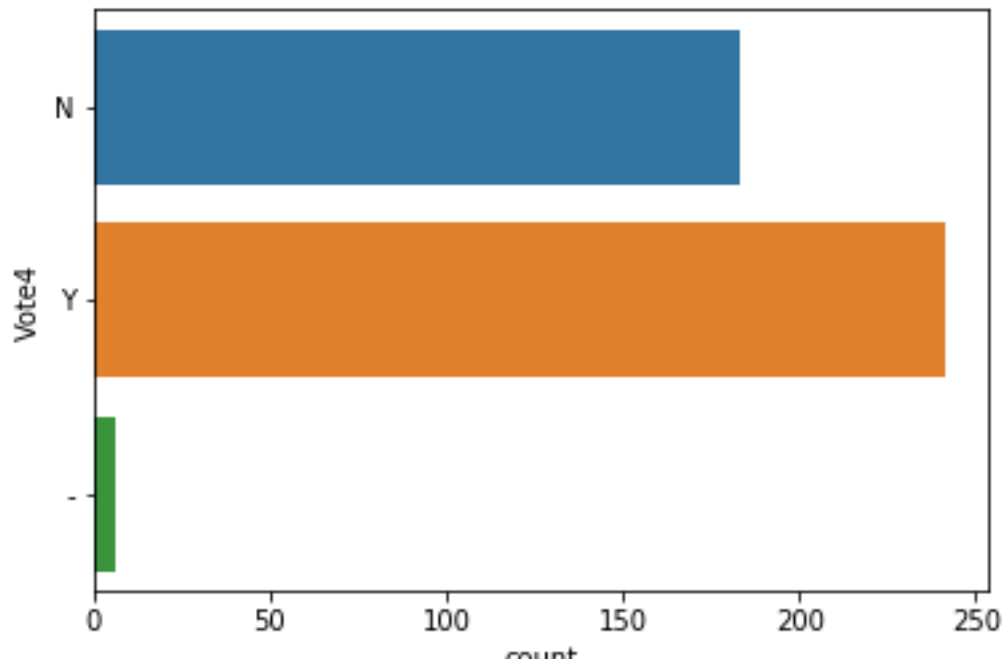




## VOTE4:

### Code:

```
sns.countplot(y="Vote4",data=vote)
```



*y Prof Chowdhury*



## DATA ANALYSIS

*Prof Chowdhury*



## ➤ NAIVE BAYES:

Naive Bayes is a simple but surprisingly powerful algorithm for predictive modeling.

In machine learning we are often interested in selecting the best hypothesis (h) given data (d).

In a classification problem, our hypothesis (h) may be the class to assign for a new data instance (d).

One of the easiest ways of selecting the most probable hypothesis given the data that we have that we can use as our prior knowledge about the problem. Bayes' Theorem provides a way that we can calculate the probability of a hypothesis given our prior knowledge.

Bayes' Theorem is stated as:

$$P(h|d) = (P(d|h) * P(h)) / P(d)$$

Where

- **P (h | d)** is the probability of hypothesis h given the data d. This is called the posterior probability.
- **P (d | h)** is the probability of data d given that the hypothesis h was true.
- **P (h)** is the probability of hypothesis h being true (regardless of the data). This is called the prior probability of h.
- **P (d)** is the probability of the data (regardless of the hypothesis).

You can see that we are interested in calculating the posterior probability of P (h | d) from the prior probability



After calculating the posterior probability for a number of different hypotheses, you can select the hypothesis with the highest probability. This is the maximum probable hypothesis and may formally be called the maximum a posteriori (MAP) hypothesis.

This can be written as:

$$\text{MAP} (h) = \max (P (h | d))$$

Or

$$\text{MAP} (h) = \max ((P (d | h) * P (h)) / P (d))$$

Or

$$\text{MAP} (h) = \max (P (d | h) * P (h))$$

The  $P(d)$  is a normalizing term which allows us to calculate the probability. We can drop it when we are interested in the most probable hypothesis as it is constant and only used to normalize.

Back to classification, if we have an even number of instances in each class in our training data, then the probability of each class (e.g.  $P (h)$ ) will be equal. Again, this would be a constant term in our equation and we could drop it so that we end up with:

$$\text{MAP} (h) = \max (P (d | h))$$

*y. Pradyumn Chowdhury*



## **Naive Bayes Classifier:**

Naive Bayes is a classification algorithm for binary (two-class) and multi-class classification problems. The technique is easiest to understand when described using binary or categorical input values.

It is called *naive Bayes* or *idiot Bayes* because the calculation of the probabilities for each hypothesis are simplified to make their calculation tractable. Rather than attempting to calculate the values of each attribute value  $P(d1, d2, d3 | h)$ , they are assumed to be conditionally independent given the target value and calculated as  $P(d1 | h) * P(d2 | H)$  and so on.

This is a very strong assumption that is most unlikely in real data, i.e. that the attributes do not interact. Nevertheless, the approach performs surprisingly well on data where this assumption does not hold.

## **Representation Used By Naive Bayes Models:**

The representation for naive Bayes is probabilities.

A list of probabilities are stored to file for a learned naive Bayes model. This includes:

- Class Probabilities: The probabilities of each class in the training dataset.
- Conditional Probabilities: The conditional probabilities of each input value given each class value.



## **Learn a Naive Bayes Model From Data:**

Learning a naive Bayes model from training data is fast.

Training is fast because only the probability of each class and the probability of each class given different input (x) values need to be calculated. No coefficients need to be fitted by optimization procedures.

## **Calculating Class Probabilities:**

The class probabilities are simply the frequency of instances that belong to each class divided by the total number of instances.

For example in a binary classification the probability of an instance belonging to class 1 would be calculated as:

$$P(\text{class}=1) = \text{count}(\text{class}=1) / (\text{count}(\text{class}=0) + \text{count}(\text{class}=1))$$

In the simplest case each class would have the probability of 0.5 or 50% for a binary classification problem with the same number of instances in each class.

## **Calculating Conditional Probabilities:**

The conditional probabilities are the frequency of each attribute value for a given class value divided by the total number of instances in that class.



For example, if a “*weather*” attribute had the values “*sunny*” and “*rainy*” and the class attribute had the class values “*go-out*” and “*stay-home*”, then the conditional probabilities of each weather value for each class value could be calculated as:

- $P(\text{weather}=\text{sunny} \mid \text{class}=\text{go-out}) = \frac{\text{count}(\text{instances with weather=sunny and class=go-out})}{\text{count}(\text{instances with class=go-out})}$
- $P(\text{weather}=\text{sunny} \mid \text{class}=\text{stay-home}) = \frac{\text{count}(\text{instances with weather=sunny and class=stay-home})}{\text{count}(\text{instances with class=stay-home})}$
- $P(\text{weather}=\text{rainy} \mid \text{class}=\text{go-out}) = \frac{\text{count}(\text{instances with weather=rainy and class=go-out})}{\text{count}(\text{instances with class=go-out})}$
- $P(\text{weather}=\text{rainy} \mid \text{class}=\text{stay-home}) = \frac{\text{count}(\text{instances with weather=rainy and class=stay-home})}{\text{count}(\text{instances with class=stay-home})}$

### **Make Predictions With a Naive Bayes Model:**

Given a naive Bayes model, you can make predictions for new data using Bayes theorem.

$$\text{MAP}(h) = \max(P(d \mid h) * P(h))$$

Using our example above, if we had a new instance with the *weather* of *sunny*, we can calculate:

$$\text{go-out} = P(\text{weather}=\text{sunny} \mid \text{class}=\text{go-out}) * P(\text{class}=\text{go-out})$$

$$\text{stay-home} = P(\text{weather}=\text{sunny} \mid \text{class}=\text{stay-home}) * P(\text{class}=\text{stay-home})$$

We can choose the class that has the largest calculated value. We can turn these values into probabilities by normalizing them as follows:



$$P(\text{go-out} | \text{weather=sunny}) = \text{go-out} / (\text{go-out} + \text{stay-home})$$
$$P(\text{stay-home} | \text{weather=sunny}) = \text{stay-home} / (\text{go-out} + \text{stay-home})$$

If we had more input variables we could extend the above example. For example, pretend we have a “*car*” attribute with the values “*working*” and “*broken*”. We can multiply this probability into the equation.

For example below is the calculation for the “go-out” class label with the addition of the car input variable set to “working”:

$$\text{go-out} = P(\text{weather=sunny} | \text{class=go-out}) * P(\text{car=working} | \text{class=go-out}) * P(\text{class=go-out})$$

### **Gaussian Naive Bayes:**

Naive Bayes can be extended to real-valued attributes, most commonly by assuming a Gaussian distribution.

This extension of naive Bayes is called Gaussian Naive Bayes. Other functions can be used to estimate the distribution of the data, but the Gaussian (or Normal distribution) is the easiest to work with because you only need to estimate the mean and the standard deviation from your training data.

*Prof Chowdhury*





### **Representation for Gaussian Naive Bayes:**

Above, we calculated the probabilities for input values for each class using a frequency. With real-valued inputs, we can calculate the mean and standard deviation of input values (x) for each class to summarize the distribution.

This means that in addition to the probabilities for each class, we must also store the mean and standard deviations for each input variable for each class.

### **Learn a Gaussian Naive Bayes Model From Data:**

This is as simple as calculating the mean and standard deviation values of each input variable (x) for each class value.

$$\text{mean}(x) = 1/n * \text{sum}(x)$$

Where n is the number of instances and x are the values for an input variable in your training data.

We can calculate the standard deviation using the following equation:

$$\text{standard deviation}(x) = \text{sqrt}(1/n * \text{sum}(xi - \text{mean}(x))^2)$$

This is the square root of the average squared difference of each value of x from the mean value of x, where n is the number of instances, sqrt() is the square root function, sum() is the sum function, xi is a specific value of the x variable for the i'th instance and mean(x) is described above, and ^2 is the square.

*Y. Pooj Chowdhury*



### Make Predictions With a Gaussian Naive Bayes Model:

Probabilities of new x values are calculated using the Gaussian Probability Density Function(PDF). When making predictions these parameters can be plugged into the Gaussian PDF with a new input for the variable, and in return the Gaussian PDF will provide an estimate of the probability of that new input value for that class.

$$\text{pdf}(x, \text{mean}, \text{sd}) = (1 / (\text{sqrt}(2 * \text{PI}) * \text{sd})) * \exp(-((x-\text{mean}^2)/(2*\text{sd}^2)))$$

Where pdf(x) is the Gaussian PDF, sqrt() is the square root, mean and sd are the mean and standard deviation calculated above, PI is the numerical constant, exp() is the numerical constant e or Euler's number raised to power and x is the input value for the input variable. We can then plug in the probabilities into the equation above to make predictions with real-valued inputs.

For example, adapting one of the above calculations with numerical values for weather and car:

$$\text{go-out} = P(\text{pdf}(\text{weather}) | \text{class}=\text{go-out}) * P(\text{pdf}(\text{car}) | \text{class}=\text{go-out}) * P(\text{class}=\text{go-out})$$

### Advantages / Disadvantages:

Naives Bayes is a stable algorithm. A small change in the the training data will not make a big change in the model.



### ➤ **Independent:**

All features are independent given the target Y.

The Naive Bayes assumption say that the probabilities of the different event (ie predictors attributes) are completely independent given the class value. (i.e., knowing the value of one attribute says nothing about the value of another if the class is known)

The attributes must be independent otherwise they will get too much influence on the accuracy. If an attribute are dependent with the target, it's the same that to determine the accuracy of only this attribute.

### ➤ **Equal:**

The event (ie predictors attributes) are equally important a priori.

If the accuracy remains exactly the same when an attribute is removed, it seems that this attribute is irrelevant to the outcome of the class (target).

### **When NAIVE BAYES is used?**

- To mark an email as spam, or not spam ?
- Classify a news article about technology, politics, or sports ?
- Check a piece of text expressing positive emotions, or negative emotions?

*Prof Chowdhury*



## ➤ K NEAREST NEIGHBORS:

“Nearest-neighbor” learning is also known as “Instance-based” learning.

K-Nearest Neighbors, or KNN, is a family of simple:

- classification
- and regression algorithms

based on Similarity (Distance) calculation between instances.

Nearest Neighbor implements rote learning. It's based on a local average calculation.

### KNN Model Representation:

The model representation for KNN is the entire training dataset.

KNN has no model other than storing the entire dataset, so there is no learning required.

Efficient implementations can store the data using complex data structures like k-d trees to make look-up and matching of new patterns during prediction efficient.

Because the entire training dataset is stored, you may want to think carefully about the consistency of your training data. It might be a good idea to curate it, update it often as new data becomes available and remove erroneous and outlier data.

*Y. Prashant Chowdhury*



## Making Predictions with KNN:

KNN makes predictions using the training dataset directly.

Predictions are made for a new instance (x) by searching through the entire training set for the K most similar instances (the neighbors) and summarizing the output variable for those K instances. For regression this might be the mean output variable, in classification this might be the mode (or most common) class value.

To determine which of the K instances in the training dataset are most similar to a new input a distance measure is used. For real-valued input variables, the most popular distance measure is Euclidean distance.

Euclidean distance is calculated as the square root of the sum of the squared differences between a new point (x) and an existing point (xi) across all input attributes j.

$$\text{EuclideanDistance}(x, x_i) = \sqrt{\sum (x_j - x_{ij})^2}$$

Euclidean is a good distance measure to use if the input variables are similar in type (e.g. all measured widths and heights). Manhattan distance is a good measure to use if the input variables are not similar in type (such as age, gender, height, etc.).

The value for K can be found by algorithm tuning. It is a good idea to try many different values for K (e.g. values 1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31, 33, 35, 37, 39, 41, 43, 45, 47, 49, 51, 53, 55, 57, 59, 61, 63, 65, 67, 69, 71, 73, 75, 77, 79, 81, 83, 85, 87, 89, 91, 93, 95, 97, 99, 101, 103, 105, 107, 109, 111, 113, 115, 117, 119, 121, 123, 125, 127, 129, 131, 133, 135, 137, 139, 141, 143, 145, 147, 149, 151, 153, 155, 157, 159, 161, 163, 165, 167, 169, 171, 173, 175, 177, 179, 181, 183, 185, 187, 189, 191, 193, 195, 197, 199, 201, 203, 205, 207, 209, 211, 213, 215, 217, 219, 221, 223, 225, 227, 229, 231, 233, 235, 237, 239, 241, 243, 245, 247, 249, 251, 253, 255, 257, 259, 261, 263, 265, 267, 269, 271, 273, 275, 277, 279, 281, 283, 285, 287, 289, 291, 293, 295, 297, 299, 301, 303, 305, 307, 309, 311, 313, 315, 317, 319, 321, 323, 325, 327, 329, 331, 333, 335, 337, 339, 341, 343, 345, 347, 349, 351, 353, 355, 357, 359, 361, 363, 365, 367, 369, 371, 373, 375, 377, 379, 381, 383, 385, 387, 389, 391, 393, 395, 397, 399, 401, 403, 405, 407, 409, 411, 413, 415, 417, 419, 421, 423, 425, 427, 429, 431, 433, 435, 437, 439, 441, 443, 445, 447, 449, 451, 453, 455, 457, 459, 461, 463, 465, 467, 469, 471, 473, 475, 477, 479, 481, 483, 485, 487, 489, 491, 493, 495, 497, 499, 501, 503, 505, 507, 509, 511, 513, 515, 517, 519, 521, 523, 525, 527, 529, 531, 533, 535, 537, 539, 541, 543, 545, 547, 549, 551, 553, 555, 557, 559, 561, 563, 565, 567, 569, 571, 573, 575, 577, 579, 581, 583, 585, 587, 589, 591, 593, 595, 597, 599, 601, 603, 605, 607, 609, 611, 613, 615, 617, 619, 621, 623, 625, 627, 629, 631, 633, 635, 637, 639, 641, 643, 645, 647, 649, 651, 653, 655, 657, 659, 661, 663, 665, 667, 669, 671, 673, 675, 677, 679, 681, 683, 685, 687, 689, 691, 693, 695, 697, 699, 701, 703, 705, 707, 709, 711, 713, 715, 717, 719, 721, 723, 725, 727, 729, 731, 733, 735, 737, 739, 741, 743, 745, 747, 749, 751, 753, 755, 757, 759, 761, 763, 765, 767, 769, 771, 773, 775, 777, 779, 781, 783, 785, 787, 789, 791, 793, 795, 797, 799, 801, 803, 805, 807, 809, 811, 813, 815, 817, 819, 821, 823, 825, 827, 829, 831, 833, 835, 837, 839, 841, 843, 845, 847, 849, 851, 853, 855, 857, 859, 861, 863, 865, 867, 869, 871, 873, 875, 877, 879, 881, 883, 885, 887, 889, 891, 893, 895, 897, 899, 901, 903, 905, 907, 909, 911, 913, 915, 917, 919, 921, 923, 925, 927, 929, 931, 933, 935, 937, 939, 941, 943, 945, 947, 949, 951, 953, 955, 957, 959, 961, 963, 965, 967, 969, 971, 973, 975, 977, 979, 981, 983, 985, 987, 989, 991, 993, 995, 997, 999, 1001, 1003, 1005, 1007, 1009, 1011, 1013, 1015, 1017, 1019, 1021, 1023, 1025, 1027, 1029, 1031, 1033, 1035, 1037, 1039, 1041, 1043, 1045, 1047, 1049, 1051, 1053, 1055, 1057, 1059, 1061, 1063, 1065, 1067, 1069, 1071, 1073, 1075, 1077, 1079, 1081, 1083, 1085, 1087, 1089, 1091, 1093, 1095, 1097, 1099, 1101, 1103, 1105, 1107, 1109, 1111, 1113, 1115, 1117, 1119, 1121, 1123, 1125, 1127, 1129, 1131, 1133, 1135, 1137, 1139, 1141, 1143, 1145, 1147, 1149, 1151, 1153, 1155, 1157, 1159, 1161, 1163, 1165, 1167, 1169, 1171, 1173, 1175, 1177, 1179, 1181, 1183, 1185, 1187, 1189, 1191, 1193, 1195, 1197, 1199, 1201, 1203, 1205, 1207, 1209, 1211, 1213, 1215, 1217, 1219, 1221, 1223, 1225, 1227, 1229, 1231, 1233, 1235, 1237, 1239, 1241, 1243, 1245, 1247, 1249, 1251, 1253, 1255, 1257, 1259, 1261, 1263, 1265, 1267, 1269, 1271, 1273, 1275, 1277, 1279, 1281, 1283, 1285, 1287, 1289, 1291, 1293, 1295, 1297, 1299, 1301, 1303, 1305, 1307, 1309, 1311, 1313, 1315, 1317, 1319, 1321, 1323, 1325, 1327, 1329, 1331, 1333, 1335, 1337, 1339, 1341, 1343, 1345, 1347, 1349, 1351, 1353, 1355, 1357, 1359, 1361, 1363, 1365, 1367, 1369, 1371, 1373, 1375, 1377, 1379, 1381, 1383, 1385, 1387, 1389, 1391, 1393, 1395, 1397, 1399, 1401, 1403, 1405, 1407, 1409, 1411, 1413, 1415, 1417, 1419, 1421, 1423, 1425, 1427, 1429, 1431, 1433, 1435, 1437, 1439, 1441, 1443, 1445, 1447, 1449, 1451, 1453, 1455, 1457, 1459, 1461, 1463, 1465, 1467, 1469, 1471, 1473, 1475, 1477, 1479, 1481, 1483, 1485, 1487, 1489, 1491, 1493, 1495, 1497, 1499, 1501, 1503, 1505, 1507, 1509, 1511, 1513, 1515, 1517, 1519, 1521, 1523, 1525, 1527, 1529, 1531, 1533, 1535, 1537, 1539, 1541, 1543, 1545, 1547, 1549, 1551, 1553, 1555, 1557, 1559, 1561, 1563, 1565, 1567, 1569, 1571, 1573, 1575, 1577, 1579, 1581, 1583, 1585, 1587, 1589, 1591, 1593, 1595, 1597, 1599, 1601, 1603, 1605, 1607, 1609, 1611, 1613, 1615, 1617, 1619, 1621, 1623, 1625, 1627, 1629, 1631, 1633, 1635, 1637, 1639, 1641, 1643, 1645, 1647, 1649, 1651, 1653, 1655, 1657, 1659, 1661, 1663, 1665, 1667, 1669, 1671, 1673, 1675, 1677, 1679, 1681, 1683, 1685, 1687, 1689, 1691, 1693, 1695, 1697, 1699, 1701, 1703, 1705, 1707, 1709, 1711, 1713, 1715, 1717, 1719, 1721, 1723, 1725, 1727, 1729, 1731, 1733, 1735, 1737, 1739, 1741, 1743, 1745, 1747, 1749, 1751, 1753, 1755, 1757, 1759, 1761, 1763, 1765, 1767, 1769, 1771, 1773, 1775, 1777, 1779, 1781, 1783, 1785, 1787, 1789, 1791, 1793, 1795, 1797, 1799, 1801, 1803, 1805, 1807, 1809, 1811, 1813, 1815, 1817, 1819, 1821, 1823, 1825, 1827, 1829, 1831, 1833, 1835, 1837, 1839, 1841, 1843, 1845, 1847, 1849, 1851, 1853, 1855, 1857, 1859, 1861, 1863, 1865, 1867, 1869, 1871, 1873, 1875, 1877, 1879, 1881, 1883, 1885, 1887, 1889, 1891, 1893, 1895, 1897, 1899, 1901, 1903, 1905, 1907, 1909, 1911, 1913, 1915, 1917, 1919, 1921, 1923, 1925, 1927, 1929, 1931, 1933, 1935, 1937, 1939, 1941, 1943, 1945, 1947, 1949, 1951, 1953, 1955, 1957, 1959, 1961, 1963, 1965, 1967, 1969, 1971, 1973, 1975, 1977, 1979, 1981, 1983, 1985, 1987, 1989, 1991, 1993, 1995, 1997, 1999, 2001, 2003, 2005, 2007, 2009, 2011, 2013, 2015, 2017, 2019, 2021, 2023, 2025, 2027, 2029, 2031, 2033, 2035, 2037, 2039, 2041, 2043, 2045, 2047, 2049, 2051, 2053, 2055, 2057, 2059, 2061, 2063, 2065, 2067, 2069, 2071, 2073, 2075, 2077, 2079, 2081, 2083, 2085, 2087, 2089, 2091, 2093, 2095, 2097, 2099, 2101, 2103, 2105, 2107, 2109, 2111, 2113, 2115, 2117, 2119, 2121, 2123, 2125, 2127, 2129, 2131, 2133, 2135, 2137, 2139, 2141, 2143, 2145, 2147, 2149, 2151, 2153, 2155, 2157, 2159, 2161, 2163, 2165, 2167, 2169, 2171, 2173, 2175, 2177, 2179, 2181, 2183, 2185, 2187, 2189, 2191, 2193, 2195, 2197, 2199, 2201, 2203, 2205, 2207, 2209, 2211, 2213, 2215, 2217, 2219, 2221, 2223, 2225, 2227, 2229, 2231, 2233, 2235, 2237, 2239, 2241, 2243, 2245, 2247, 2249, 2251, 2253, 2255, 2257, 2259, 2261, 2263, 2265, 2267, 2269, 2271, 2273, 2275, 2277, 2279, 2281, 2283, 2285, 2287, 2289, 2291, 2293, 2295, 2297, 2299, 2301, 2303, 2305, 2307, 2309, 2311, 2313, 2315, 2317, 2319, 2321, 2323, 2325, 2327, 2329, 2331, 2333, 2335, 2337, 2339, 2341, 2343, 2345, 2347, 2349, 2351, 2353, 2355, 2357, 2359, 2361, 2363, 2365, 2367, 2369, 2371, 2373, 2375, 2377, 2379, 2381, 2383, 2385, 2387, 2389, 2391, 2393, 2395, 2397, 2399, 2401, 2403, 2405, 2407, 2409, 2411, 2413, 2415, 2417, 2419, 2421, 2423, 2425, 2427, 2429, 2431, 2433, 2435, 2437, 2439, 2441, 2443, 2445, 2447, 2449, 2451, 2453, 2455, 2457, 2459, 2461, 2463, 2465, 2467, 2469, 2471, 2473, 2475, 2477, 2479, 2481, 2483, 2485, 2487, 2489, 2491, 2493, 2495, 2497, 2499, 2501, 2503, 2505, 2507, 2509, 2511, 2513, 2515, 2517, 2519, 2521, 2523, 2525, 2527, 2529, 2531, 2533, 2535, 2537, 2539, 2541, 2543, 2545, 2547, 2549, 2551, 2553, 2555, 2557, 2559, 2561, 2563, 2565, 2567, 2569, 2571, 2573, 2575, 2577, 2579, 2581, 2583, 2585, 2587, 2589, 2591, 2593, 2595, 2597, 2599, 2601, 2603, 2605, 2607, 2609, 2611, 2613, 2615, 2617, 2619, 2621, 2623, 2625, 2627, 2629, 2631, 2633, 2635, 2637, 2639, 2641, 2643, 2645, 2647, 2649, 2651, 2653, 2655, 2657, 2659, 2661, 2663, 2665, 2667, 2669, 2671, 2673, 2675, 2677, 2679, 2681, 2683, 2685, 2687, 2689, 2691, 2693, 2695, 2697, 2699, 2701, 2703, 2705, 2707, 2709, 2711, 2713, 2715, 2717, 2719, 2721, 2723, 2725, 2727, 2729, 2731, 2733, 2735, 2737, 2739, 2741, 2743, 2745, 2747, 2749, 2751, 2753, 2755, 2757, 2759, 2761, 2763, 2765, 2767, 2769, 2771, 2773, 2775, 2777, 2779, 2781, 2783, 2785, 2787, 2789, 2791, 2793, 2795, 2797, 2799, 2801, 2803, 2805, 2807, 2809, 2811, 2813, 2815, 2817, 2819, 2821, 2823, 2825, 2827, 2829, 2831, 2833, 2835, 2837, 2839, 2841, 2843, 2845, 2847, 2849, 2851, 2853, 2855, 2857, 2859, 2861, 2863, 2865, 2867, 2869, 2871, 2873, 2875, 2877, 2879, 2881, 2883, 2885, 2887, 2889, 2891, 2893, 2895, 2897, 2899, 2901, 2903, 2905, 2907, 2909, 2911, 2913, 2915, 2917, 2919, 2921, 2923, 2925, 2927, 2929, 2931, 2933, 2935, 2937, 2939, 2941, 2943, 2945, 2947, 2949, 2951, 2953, 2955, 2957, 2959, 2961, 2963, 2965, 2967, 2969, 2971, 2973, 2975, 2977, 2979, 2981, 2983, 2985, 2987, 2989, 2991, 2993, 2995, 2997, 2999, 3001, 3003, 3005, 3007, 3009, 3011, 3013, 3015, 3017, 3019, 3021, 3023, 3025, 3027, 3029, 3031, 3033, 3035, 3037, 3039, 3041, 3043, 3045, 3047, 3049, 3051, 3053, 3055, 3057, 3059, 3061, 3063, 3065, 3067, 3069, 3071, 3073, 3075, 3077, 3079, 3081, 3083, 3085, 3087, 3089, 3091, 3093, 3095, 3097, 3099, 3101, 3103, 3105, 3107, 3109, 3111, 3113, 3115, 3117, 3119, 3121, 3123, 3125, 3127, 3129, 3131, 3133, 3135, 3137, 3139, 3141, 3143, 3145, 3147, 3149, 3151, 3153, 3155, 3157, 3159, 3161, 3163, 3165, 3167, 3169, 3171, 3173, 3175, 3177, 3179, 3181, 3183, 3185, 3187, 3189, 3191, 3193, 3195, 3197, 3199, 3201, 3203, 3205, 3207, 3209, 3211, 3213, 3215, 3217, 3219, 3221, 3223, 3225, 3227, 3229, 3231, 3233, 3235, 3237, 3239, 3241, 3243, 3245, 3247, 3249, 3251, 3253, 3255, 3257, 3259, 3261, 3263, 3265, 3267, 3269, 3271, 3273, 3275, 3277, 3279, 3281, 3283, 3285, 3287, 3289, 3291, 3293, 3295, 3297, 3299, 3301, 3303, 3305, 3307, 3309, 3311, 3313, 3315, 3317, 3319, 3321, 3323, 3325, 3327, 3329, 3331, 3333, 3335, 3337, 3339, 3341, 3343, 3345, 3347, 3349, 3351, 3353, 3355, 3357, 3359, 3361, 3363, 3365, 3367, 3369, 3371, 3373, 3375, 3377, 3379, 3381, 3383, 3385, 3387, 3389, 3391, 3393, 3395, 3397, 3399, 3401, 3403, 3405, 3407, 3409, 3411, 3413, 3415, 3417, 3419, 3421, 3423, 3425, 3427, 3429, 3431, 3433, 3435, 3437, 3439, 3441, 3443, 3445, 3447, 3449, 3451, 3453, 3455, 3457, 3459, 3461, 3463, 3465, 3467, 3469, 3471, 3473, 3475, 3477, 3479, 3481, 3483, 3485, 3487, 3489, 3491, 3493, 3495, 3497, 3499, 3501, 3503, 3505, 3507, 3509, 3511, 3513, 3515, 3517, 3519, 3521, 3523, 3525, 3527, 3529, 3531, 3533, 3535, 3537, 3539, 3541, 3543, 3545, 3547, 3549, 3551, 3553, 3555, 3557, 3559, 3561, 3563, 3565, 3567, 3569, 3571, 3573, 3575, 3577, 3579, 3581, 3583, 3585, 3587, 3589, 3591, 3593, 3595, 3597, 3599, 3601, 3603, 3605, 3607, 3609, 3611, 3613, 3615, 3617, 3619, 3621, 3623, 3625, 3627, 3629, 3631, 3633, 3635, 3637, 3639, 3641, 3643, 3645, 3647, 3649, 3651, 3653, 3655, 3657, 3659, 3661, 3663, 3665, 3667, 3669, 3671, 3673, 3675, 3677, 3679, 3681, 3683, 3685, 3687, 3689, 3691, 3693, 3695, 3697, 3699, 3701, 3703, 3705, 3707, 3709, 3711, 3713, 3715, 3717, 3719, 3721, 3723, 3725, 3727, 3729, 3731, 3733, 3735, 3737, 3739, 3741, 3743, 3745, 3747, 3749, 3751, 3753, 3755, 3757, 3759, 3761, 3763, 3765, 3767, 3769, 3771, 3773, 3775, 3777, 3779, 3781, 3783, 3785, 3787, 3789, 3791, 3793, 3795, 3797, 3799, 3801, 3803, 3805, 3807, 3809, 3811, 3813, 3815, 3817, 3819, 3821, 3823, 3825, 3827, 3829, 3831, 3833, 3835, 3837, 3839, 3841, 3843, 3845, 3847, 3849, 3851, 3853, 3855, 3857, 3859, 3861, 3863, 3865, 3867, 3869, 3871, 3873, 3875, 3877, 3879, 3881, 3883, 3885, 3887, 3889, 3891, 3893, 3895, 3897, 3899, 3901, 3903, 3905, 3907, 3909, 3911, 3913, 3915, 3917, 3919, 3921, 3923, 3925, 3927, 3929, 3931, 3933, 3935, 3937, 3939, 3941, 3943, 3945, 3947, 3949, 3951, 3953, 3955, 3957, 3959, 3961, 3963, 3965, 3967, 3969, 3971, 3973, 3975, 3977, 3979, 3981, 3983, 3985, 3987, 3989, 3991, 3993, 3995, 3997, 3999, 4001, 4003, 4005, 4007, 4009, 4011, 4013, 4015, 4017, 4019, 4021, 4023, 4025, 4027, 4029, 4031, 4033, 4035, 4037, 4039, 4041, 4043, 4045, 4047, 4049, 4051, 4053, 4055, 4057, 4059, 4061, 4063, 4065, 4067, 4069, 4071, 4073, 4075, 4077, 4079, 4081, 4083, 4085, 4087, 4089, 4091, 4093, 4095, 4097, 4099, 4101, 4103, 4105, 4107, 4109, 4111, 4113, 4115, 4117, 4119, 4121, 4123, 4125, 4127, 4129, 4131, 4133, 4135, 4137,

The computational complexity of KNN increases with the size of the training dataset. For very large training sets, KNN can be made stochastic by taking a sample from the training dataset from which to calculate the K-most similar instances.

KNN has been around for a long time and has been very well studied. As such, different disciplines have different names for it, for example:

- **Instance-Based Learning:** The raw training instances are used to make predictions. As such KNN is often referred to as instance-based learning or a case-based learning (where each training instance is a case from the problem domain).
  - **Lazy Learning:** No learning of the model is required and all of the work happens at the time a prediction is requested. As such, KNN is often referred to as a lazy learning algorithm.
  - **Non-Parametric:** KNN makes no assumptions about the functional form of the problem being solved. As such KNN is referred to as a non-parametric machine learning algorithm.
- KNN can be used for regression and classification problems.

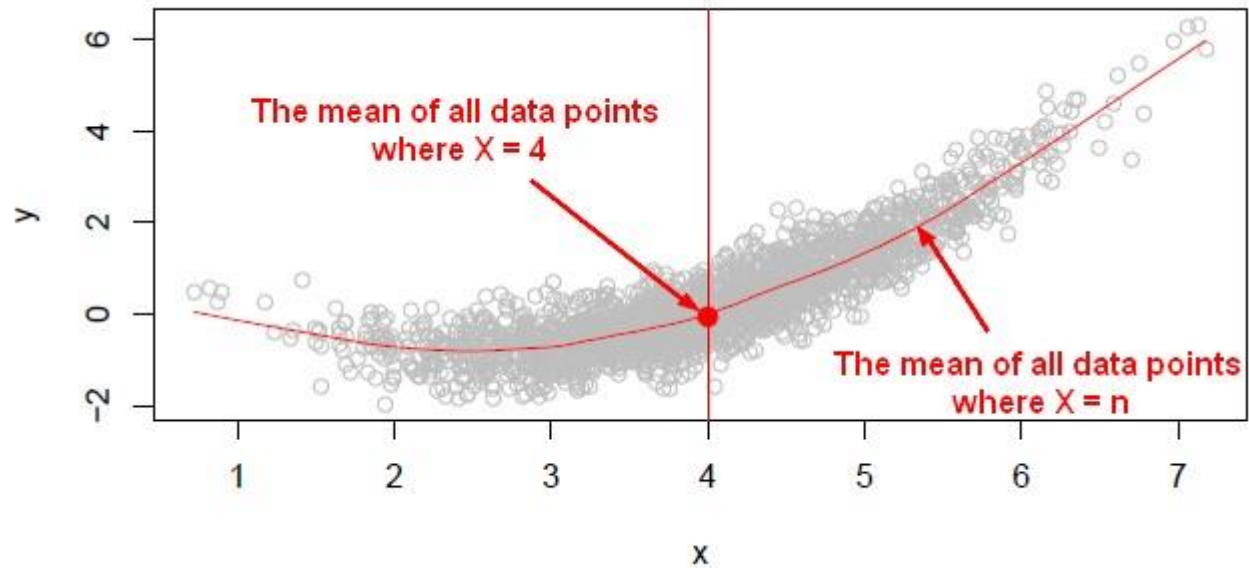
### KNN for Regression:

When KNN is used for regression problems the prediction is based on the mean or the median of the K-most similar instances.

*Y. Pooj Chowdhury*



➤ **Mean:**



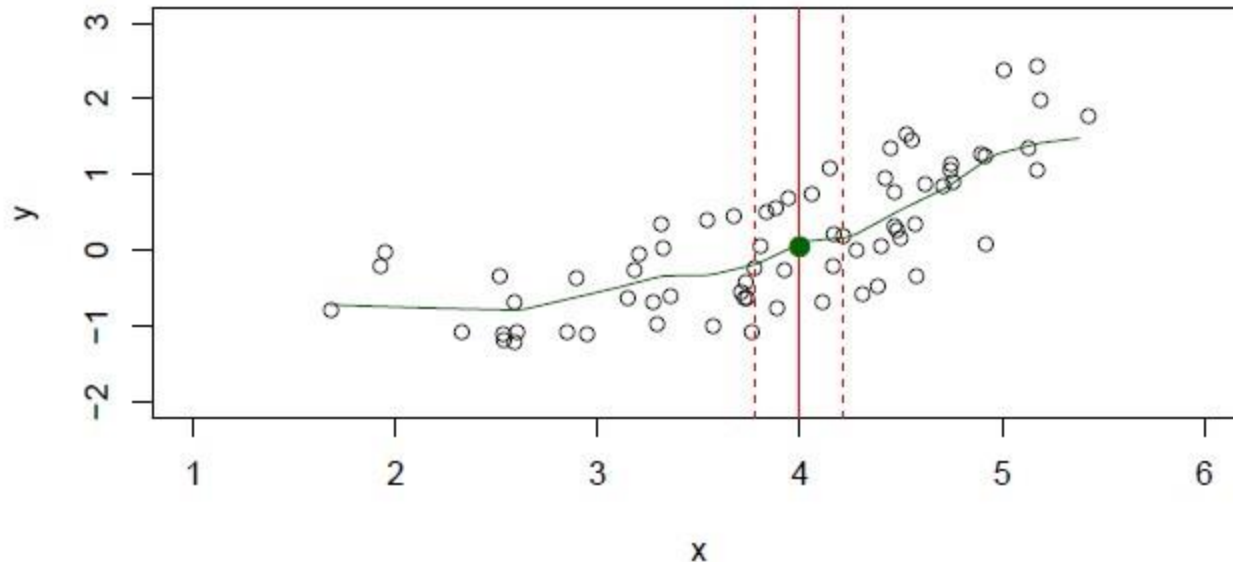
Of all Y value for a X.

➤ **Neighborhood:**

When they are too few data points for a X value in order to calculate the mean, the neighbourhood can be used.

*Prof Chowdhury*





Nearest neighbour averaging can be pretty good for a small number of variable (See accurate ) and large N(in order to get enough point to calculate the average).

➤ **KNN for Classification:**

When KNN is used for classification, the output can be calculated as the class with the highest frequency from the K-most similar instances. Each instance in essence votes for their class and the class with the most votes is taken as the prediction.

Class probabilities can be calculated as the normalized frequency of samples that belong to each class in the set of K most similar instances for a new data instance. For example, in a binary classification

*y. Prashant Chowdhury*

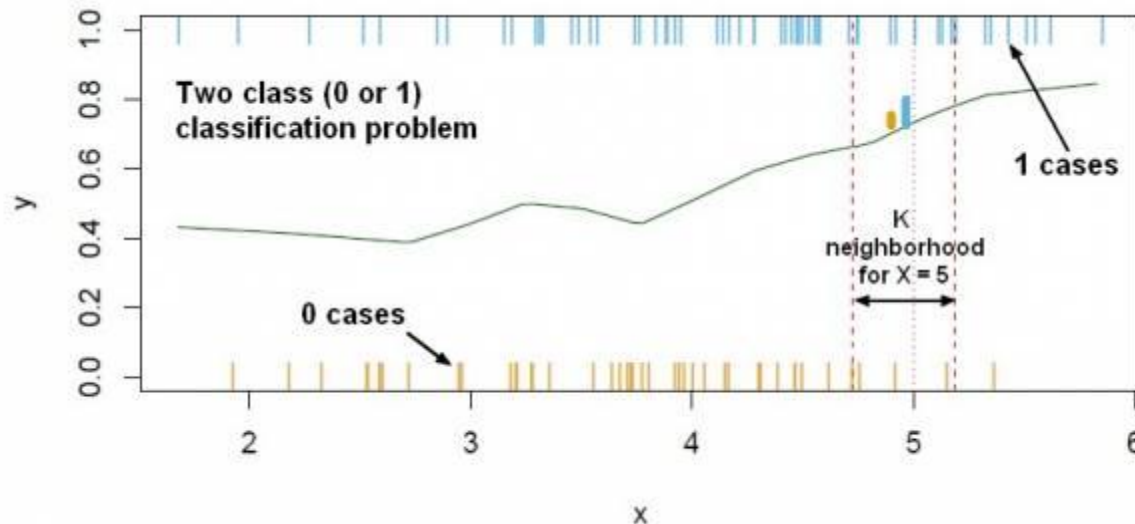




$$p(\text{class}=0) = \text{count}(\text{class}=0) / (\text{count}(\text{class}=0) + \text{count}(\text{class}=1))$$

If you are using K and you have an even number of classes (e.g. 2) it is a good idea to choose a K value with an odd number to avoid a tie. And the inverse, use an even number for K when you have an odd number of classes.

Ties can be broken consistently by expanding K by 1 and looking at the class of the next most similar instance in the training dataset.



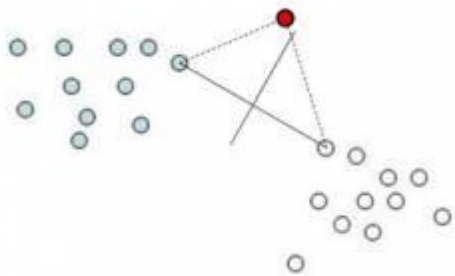
### ➤ Decision boundary:

The nearest neighbor method produces a linear decision boundary. It's a little bit more complicated as it produces a piece-wise linear decision of the decision boundary with sometimes a bunch of little linear segments.

*Prof Chowdhury*



The perpendicular bisector of the line that joins the two closest points:



➤ **Noisy instances:**

Noisy instance are instance with a bad target class.

If the dataset is noisy , then by accident we might find an incorrectly classified training instance as the nearest one to our test instance.

You can guard against that by using:

- Majority vote over K nearest neighbours instances (k might be 3 or 4). You look for the 3 or 5 nearest neighbours and choose the majority class amongst those when classifying an unknown point. See K parameter
- Weight instances according to prediction accuracy
- Identification of reliable “prototypes” for each class

**K Parameter:**

The K param

*y. Prasad Chowdhury*

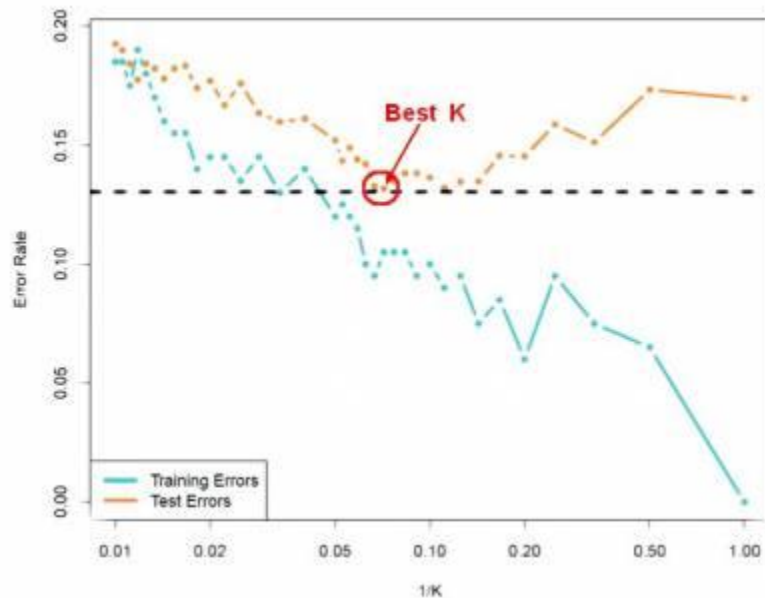
method against noisy dataset.

### ➤ Best k value:

An obvious issue with k nearest neighbour is how to choose a suitable value for the number of nearest neighbours used.

- If it's too small, the method is susceptible to noise in the data.
- If it's too large, the decision is smeared out, covering too great an area of instance space

Weka uses cross-validation to select the best value.



The X axis shows the formula 1/K because since K is the number of neighbours,  $1/k \leftarrow 1$ .

*Prof Chowdhury*



➤ **With a K close to the size of the whole data set:**

If we set  $k$  to be an extreme value, close to the size of the whole data set, then we're taking the distance of the test instance to all of points in the dataset and averaging those which will probably gives us something close to the baseline accuracy.

**Advantage/Inconvenient:**

---

➤ **Curse of Dimensionality:**

KNN works well with a small number of input variables ( $p$ ), but struggles when the number of inputs is very large.

Each input variable can be considered a dimension of a  $p$ -dimensional input space. For example, if you had two input variables  $x_1$  and  $x_2$ , the input space would be 2-dimensional.

As the number of dimensions increases the volume of the input space increases at an exponential rate.

In high dimensions, points that may be similar may have very large distances. All points will be far away from each other and our intuition for distances in simple 2 and 3-dimensional spaces breaks down. This might feel unintuitive at first, but this general problem is called the “Curse of Dimensionality”.

*y. Prashant Chowdhury*



➤ **Slow:**

slow as you need to scan entire training data to make each prediction.

➤ **Checking wheather the data is noisy:**

By changing the k data set if the accuracy changes a lot between 1 and 5 k it's may be a noisy data set If the data set is noisy, the accuracy figures improves as k got little bit larger but then it would be starting to decrease again.

## **When KNN is used?**

- k-NN is often used in search applications to look for “similar” items; that is, when task is some form of “find items similar to this one”.
- The biggest use case of k-NN search might be Recommender Systems. If a user likes a particular item, then k-NN can recommend similar items for them. To find similar items, the set of users who like each item are compared , if a similar set of users like two different items, then the items themselves are probably similar!

*Y. Prashant Chowdhury*



CODE

*Prof Chowdhury*



```
import pandas as pd
import os
import re
import seaborn as sns
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn import preprocessing
from sklearn import neighbors

os.chdir ( "c:/datasets" )
reg=re.compile(r'\s\s+')
to_file=open("formatted_vote.csv",mode="w",encoding="utf-8")
with open("vote.csv",encoding="utf-8") as a_file:
    for a_line in a_file:
        formatted_rec=re.sub( r'\s\s+', ', ', a_line )
        to_file.write(formatted_rec)
to_file.close()
```

  
  
Document sign date :Apr 23, 2018

```
#we have to manually edit the column names
vote=pd.read_csv("formatted_vote.csv")
vote=vote[(vote.Party != "Y") & (vote.Party != "N")]

sns.countplot(y="Party",data=vote)
sns.countplot(y="Vote1",data=vote)
sns.countplot(y="Vote2",data=vote)
sns.countplot(y="Vote3",data=vote)
sns.countplot(y="Vote4",data=vote)

vote=vote.dropna(axis=0,how='any',inplace=False)
```

#####

*y. Prashant Chowdhury*





#####

```
def_r=(vote.Party == "R")  
def_v1=(vote.Vote1 == "Y")  
def_v2=(vote.Vote2 == "Y")  
def_v3=(vote.Vote3 == "Y")  
def_v4=(vote.Vote4 == "Y")
```

```
le = preprocessing.LabelEncoder()  
bin_r=le.fit_transform(def_r)  
bin_v1=le.fit_transform(def_v1)  
bin_v2=le.fit_transform(def_v2)  
bin_v3=le.fit_transform(def_v3)  
bin_v4=le.fit_transform(def_v4)
```

*y. Prasad Chowdhury*



```
df_r = pd.DataFrame(bin_r)
```

```
df_v1 = pd.DataFrame(bin_v1)
```

```
df_v2 = pd.DataFrame(bin_v2)
```

```
df_v3 = pd.DataFrame(bin_v3)
```

```
df_v4 = pd.DataFrame(bin_v4)
```

```
df1 = df_r.assign(df_v1=df_v1.values)
```

```
df1 = df1.assign(df_v2=df_v2.values)
```

```
df1 = df1.assign(df_v3=df_v3.values)
```

```
df1 = df1.assign(df_v4=df_v4.values)
```

```
df1.columns=['Party','Vote1','Vote2','Vote3','Vote4']
```

```
Xnb=df1[["Vote1","Vote2","Vote3","Vote4"]]
```

```
ynb=df1["Party"]
```

```
nb = GaussianNB()
```

```
Xtrain, Xtest, ytrain, ytest = train_test_split(Xnb, ynb, test_size=.25, random_state=42)
```

```
nb.fit(Xtrain,ytrain)
```

```
performance_nb_test=nb.score(Xtest,ytest)
```

```
performance_nb_train=nb.score(Xtrain,ytrain)
```

```
print("")
```

```
print("Accuracy over train set in Naive Bayes: ",performance_nb_train)
```

```
print("Accuracy over test set in Naive Bayes: ",performance_nb_test)
```

```
print("")
```

```
#####
```

```
#knn accuracy
```

```
#####
```

```
vote=vote[vote.Vote1 != "-"]
```

```
vote=vote[vote.Vote2 != "-"]
```

```
vote=vote[vote.Vote3 != "-"]
```

```
vote=vote[vote.Vote4 != "-"]
```

*y Pro Chowdhury*



```
def_r=(vote.Party == "R")  
def_v1=(vote.Vote1 == "Y")  
def_v2=(vote.Vote2 == "Y")  
def_v3=(vote.Vote3 == "Y")  
def_v4=(vote.Vote4 == "Y")
```

```
le = preprocessing.LabelEncoder()  
bin_r=le.fit_transform(def_r)  
bin_v1=le.fit_transform(def_v1)  
bin_v2=le.fit_transform(def_v2)  
bin_v3=le.fit_transform(def_v3)  
bin_v4=le.fit_transform(def_v4)
```

```
df_r = pd.DataFrame(bin_r)
```

```
df_v1 = pd.L
```

*Prof Chowdhury*



```
df_v2 = pd.DataFrame(bin_v2)
```

```
df_v3 = pd.DataFrame(bin_v3)
```

```
df_v4 = pd.DataFrame(bin_v4)
```

```
df1 = df_r.assign(df_v1=df_v1.values)
```

```
df1 = df1.assign(df_v2=df_v2.values)
```

```
df1 = df1.assign(df_v3=df_v3.values)
```

```
df1 = df1.assign(df_v4=df_v4.values)
```

```
df1.columns=['Party','Vote1','Vote2','Vote3','Vote4']
```

```
X=df1[["Party","Vote1","Vote2","Vote3"]]
```

```
y=df1["Vote4"]
```

```
Xtrain, Xtest, ytrain, ytest=train_test_split(X,y,test_size=.20,random_state=42)
```

```
knn= neighbors.KNeighborsClassifier()
```

```
knn.fit(Xtrain
```

*y Prof Chowdhury*



```
performance_knn_test=knn.score(Xtest,ytest)
performance_knn_train=knn.score(Xtrain,ytrain)

print("Accuracy over train set in KNN: ",performance_knn_train)
print("Accuracy over test set in KNN: ",performance_knn_test)

print("")

#####

#predict party naive bayes

#####

os.chdir("c:/datasets")
vote=pd.read_csv("formatted_vote.csv")
vote=vote[(vote.Party != "Y") & (vote.Party != "N")]
```

*y. Prashant Chowdhury*



```
vote=vote.dropna(axis=0,how='any',inplace=False)
```

```
x={'v1y':[0,0],'v1n':[0,0],'v2y':[0,0],'v2n':[0,0],'v3y':[0,0],'v3n':[0,0],'v4y':[0,0],'v4n':[0,0],'all':[0,0]}
```

```
index=["R","D"]
```

```
freqdf=pd.DataFrame(x,index)
```

```
vote1=vote[["Party","Vote1","Vote2","Vote3","Vote4"]]
```

```
for rows in vote1.iterrows():
```

```
    if rows[1][0]=="R":
```

```
        freqdf.loc[rows[1][0],"all"]=freqdf.loc[rows[1][0],"all"]+1
```

```
    if rows[1][0]=="D":
```

```
        freqdf.loc[rows[1][0],"all"]=freqdf.loc[rows[1][0],"all"]+1
```

```
    if rows[1][1]=="Y":
```

```
        freqdf.loc[rows[1][0],"v1y"]=freqdf.loc[rows[1][0],"v1y"]+1
```

```
    if rows[1]
```

*y. Prashant Chowdhury*



```

freqdf.loc[rows[1][0],"v1n"]=freqdf.loc[rows[1][0],"v1n"]+1
if rows[1][2]=="Y":
    freqdf.loc[rows[1][0],"v2y"]=freqdf.loc[rows[1][0],"v2y"]+1
if rows[1][2]=="N":
    freqdf.loc[rows[1][0],"v2n"]=freqdf.loc[rows[1][0],"v2n"]+1
if rows[1][3]=="Y":
    freqdf.loc[rows[1][0],"v3y"]=freqdf.loc[rows[1][0],"v3y"]+1
if rows[1][3]=="N":
    freqdf.loc[rows[1][0],"v3n"]=freqdf.loc[rows[1][0],"v3n"]+1
if rows[1][4]=="Y":
    freqdf.loc[rows[1][0],"v4y"]=freqdf.loc[rows[1][0],"v4y"]+1
if rows[1][4]=="N":
    freqdf.loc[rows[1][0],"v4n"]=freqdf.loc[rows[1][0],"v4n"]+1

```

```

x={'v1y':[0,0],'v1n':[0,0],'v2y':[0,0],'v2n':[0,0],'v3y':[0,0],'v3n':[0,0],'v4y':[0,0],'v4n':[0,0],'all':[0,0]}

```

```

index=["R",

```

*Prof Chowdhury*





```

probd f=pd.DataFrame(x,index)
r=0
for row in freqdf.iterrows():
    i=1
    while i<=8:
        probdf.iloc[r,i]=row[1][i]/row[1][0]
        i=i+1
    r=r+1
probd f.iloc[0,0]=freqdf.iloc[0,0]/(freqdf.iloc[0,0]+freqdf.iloc[1,0])
probd f.iloc[1,0]=freqdf.iloc[1,0]/(freqdf.iloc[0,0]+freqdf.iloc[1,0])

print("In Naive Bayes, party prediction for last four Congressmen :")

vote2=vote.tail(4).copy()
count_yes=0
count_no=0
for row in vc

```

*y. Prashant Chowdhury*



```
v1=row[1][3]
```

```
v2=row[1][4]
```

```
v3=row[1][5]
```

```
v4=row[1][6]
```

```
prod_r=probd.f.loc["R","all"]
```

```
prod_d=probd.f.loc["D","all"]
```

```
if v1=="Y":
```

```
    prod_r=prod_r*probd.f.loc["R","v1y"]
```

```
    prod_d=prod_d*probd.f.loc["D","v1y"]
```

```
else:
```

```
    prod_r=prod_r*probd.f.loc["R","v1n"]
```

```
    prod_d=prod_d*probd.f.loc["D","v1n"]
```

```
if v2=="Y":
```

```
    prod_r=prod_r*probd.f.loc["R","v2y"]
```

```
    prod_d:
```

*y. Prasad Chowdhury*



else:

prod\_r=prod\_r\*probd.f.loc["R","v2n"]

prod\_d=prod\_d\*probd.f.loc["D","v2n"]

if v3=="Y":

prod\_r=prod\_r\*probd.f.loc["R","v3y"]

prod\_d=prod\_d\*probd.f.loc["D","v3y"]

else:

prod\_r=prod\_r\*probd.f.loc["R","v3n"]

prod\_d=prod\_d\*probd.f.loc["D","v3n"]

if v4=="Y":

prod\_r=prod\_r\*probd.f.loc["R","v4y"]

prod\_d=prod\_d\*probd.f.loc["D","v4y"]

else:

prod\_r=prod\_r\*probd.f.loc["R","v4n"]

prod\_d:

*y. Prasad Chowdhury*



```
predicted = "R" if prod_r >= prod_d else "D"
if predicted==row[1][2]:
    count_yes=count_yes+1
else:
    count_no=count_no+1
print("Predicted : ",predicted," Original : ",row[1][2])
```

```
accuracy_test=count_yes/(count_yes+count_no)
print("Accuracy is: ",accuracy_test)
```

```
print("")
```

```
#####
```

```
#knn_Predict_Vote4
```

```
#####
```

*y. Prashant Chowdhury*



```
vote2=df1.tail(4)
```

```
vote2=vote2[["Party","Vote1","Vote2","Vote3"]]
```

```
result2=knn.predict(vote2)
```

```
yresult2=df1.tail(4)
```

```
yresult2=yresult2["Vote4"]
```

```
accuracy=knn.score(vote2,yresult2)
```

```
result21= pd.DataFrame(result2)
```

```
df=result21.assign(yresult2=yresult2.values)
```

```
df.columns=['Predicted','Original']
```

```
df2=df
```

```
df2 = df2.replace( 1, "Y")
```

```
df2 = df2.replace( 0, "N")
```

*y. Prashant Chowdhury*



```
print("In knn, The last four Congressmen Vote4 will be :")
```

```
print(df2)
```

```
print("Accuracy is: ",accuracy)
```

**#End\_of\_code**

*y. Prof Chowdhury*



## RESULT:

Accuracy over train set in Naive Bayes: 0.850931677019

Accuracy over test set in Naive Bayes: 0.935185185185

Accuracy over train set in KNN: 0.708860759494

Accuracy over test set in KNN: 0.75

In Naive Bayes, party prediction for last four Congressmen:

Predicted: D Original: D

Predicted: D Original: D

Predicted: R Original: R

Predicted: D Original: R

Accuracy is: 0.75

In knn, the last four Congressmen Vote4 will be:

Predicted C

*Y. Prashant Chowdhury*



0	Y	Y
1	Y	Y
2	N	N
3	N	Y

Accuracy is: 0.75

➤ **FINAL COLUMN EXCLUSION:**

No column has been eliminated.

➤ **EXPERIMENT:**

The accuracy in Naïve Bayes over test set is 0.93 while in KNN, it is 0.75.

So, here we can see Naïve Bayes has performed a better result.

*y. Prashant Chowdhury*





## FUTURE IMPROVEMENT

For future improvement of the model & Betterment of the scores ; Deep Learning , Nural Network & other advance methods may be introduced to it.

*Y. P. Chowdhury*



## CERTIFICATE

This is to certify that **Mr. SPANDAN MAITY** of **College of Engineering and Management, Kolaghat**, registration number: **151070110050 OF 2015-2016**, has successfully completed a project on **“Machine Learning using Python”** under the guidance of **Mr. Titas Roy Chowdhury**.

---

[Mr. Titas Roy Chowdhury]



## CERTIFICATE

This is to certify that **Mr. DEBANJAN DAS** of **College of Engineering and Management, Kolaghat**, registration number: **151070110021 OF 2015-2016**, has successfully completed a project on **“Machine Learning using Python”** under the guidance of **Mr. Titas Roy Chowdhury**.

---

[Mr. Titas Roy Chowdhury]

