

ML Notes without the BS

vivian sedov

CONTENTS

I	Types of Learning	2
I-A	Supervised Learning	2
II	Introduction To NN	4
II-A	How does it work ?	4
II-A1	Step by step guide on how NN works	4
II-B	KNN	4
III	Conformal Prediction	6
III-A	Explanation of Algo	6
III-A1	Special Cases . .	7
III-A2	Assumptions to- wards Conformity	7
III-B	Conformity Measure for KNN (K=1)	8
III-C	Ranking Algo	8
III-C1	Example	8
III-D	Example	9

I. TYPES OF LEARNING

- Supervised Learning
 - Algo is trained on labeled data
 - meaning that the data has been labeled with the correct output
 - *Goal:* is to make predictions about new unseen data based on patterns learned from the training data.
 - Examples
 - * Linear Regression
 - * Logistic Regression
 - * Support vector Machines
 - * Classification problems : Binary or Multi Class
- Unsupervised Learning
 - The algo is not given any labeled data and must find the patterns and relationships on its own.
 - *Goal:* Is to discover structure in the data and to identify relationships based on that structure.
 - Examples:
 - * Clustering
 - * Dimensionality reduction [Really cool concept]
- Reinforcement Learning
 - The learning algo is not given any labeled data and must find patterns and relationships in the data on its own.
 - *Goal :* is to discover structure, but based on its own ability to generalise and understand the current scope of what it can learn from.
 - Examples:
 - * Self driving cars
 - * Those small little hoovers that move around on their own. Question your self how do they map everything.

A. Supervised Learning

Due to this module purely focusing on this, i will pick certain info that i deem viable

a) Features:

- Components of samples are features
- Each feature can be discrete or continuous $[0, \infty)$

b) Protocols used in supervised Learning:

There are two core protocols that are used within supervised learning: Batch and Online | Just google this, the explanation on the slides are bad.

Batch Learning In batch learning, there are two core steps, and stages that you have to work with

- Training: Exploration stage - You analyse the data, and the training set - and find some viable explanation for why it works
- Exploitation stage: Pretty much saying does this hypothesis work with respect to what the training data had presented.

Definition I.1 – Induction

Induction is a form of learning that involves making generalizations based on specific examples or observations. In induction, the goal is to learn a general rule or model that can be applied to new, unseen data. This is the approach used in many supervised learning algorithms, where the algorithm is trained on a labeled dataset and then makes predictions about new, unseen data based on the patterns learned from the training data.

has a model

Definition I.2 – Transduction

Transduction is a form of learning that involves making predictions about a specific instance based on the available information. In transduction, the goal is to make a prediction about a specific case rather than learning a general rule that can be applied to new data. This is the approach used in many unsupervised learning algorithms, where the algorithm is given a dataset but is not told what the correct output should be. Instead, the algorithm must find patterns and relationships in the data on its own and use that information to make a prediction about a specific instance.

No model

The main difference between induction and transduction is that induction involves learning a gen-

eral rule or model that can be applied to new data, while transduction involves making a prediction about a specific instance based on the available information.

II. INTRODUCTION TO NN

Definition II.1 – Nearest Neighbor Algorithm The nearest neighbor algorithm is a method for classification that assigns a new data point to the class of the nearest training data point. This algorithm is based on the idea that similar data points are likely to belong to the same class.

Within binary classification problems, the label space is often taken to be 0 or 1.

A. How does it work ?

Here is a mathematical example of how the nearest neighbor algorithm works: Suppose we have a training set of data points with known classes, represented as a matrix $X \in \mathbb{R}^{n \times m}$, where each row corresponds to a data point and each column corresponds to a feature. The class labels for the data points are stored in a vector $y \in \mathbb{R}^n$.

Now, suppose we have a new data point, represented as a row vector $x \in \mathbb{R}^m$, that we want to classify. To classify this point using the nearest neighbor algorithm, we need to find the data point in the training set that is closest to x . We can measure the distance between x and each data point in the training set using a distance metric, such as the Euclidean distance:

$$d(x, X_i) = \sqrt{\sum_{j=1}^m (x_j - X_{i,j})^2}$$

where X_i is the i -th row of the matrix X , and x_j and $X_{i,j}$ are the j -th elements of the vectors x and X_i , respectively.

Once we have computed the distances between x and all the data points in the training set, we can find the nearest neighbor by selecting the data point with the smallest distance:

$$\hat{y} = y_i, \text{ where } i = \operatorname{argmin}_{j=1}^n d(x, X_j)$$

In this equation, \hat{y} is the predicted class for the new data point x , and y_i is the class of the nearest neighbor.

1) Step by step guide on how NN works:

Suppose we have the following training set of data points, where each row corresponds to a data point and each column corresponds to a feature:

$$X = \begin{bmatrix} 1 & 1 & 2 & 2 & 3 & 3 & 4 & 4 \end{bmatrix}$$

The class labels for the data points are stored in the following vector:

$$y = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}$$

Now, suppose we have a new data point that we want to classify, represented as the following row vector:

$$x = \begin{bmatrix} 1.5 & 1.5 \end{bmatrix}$$

To classify this point using the nearest neighbor algorithm, we need to find the data point in the training set that is closest to x . We can measure the distance between x and each data point in the training set using the Euclidean distance, as shown in the previous equation.

For example, the distance between x and the first data point in the training set is:

$$\begin{aligned} d(x, X_1) &= \sqrt{(1.5 - 1)^2 + (1.5 - 1)^2} \\ &= \sqrt{0.5^2 + 0.5^2} = 0.70715 \end{aligned}$$

We can compute the distances between x and all the other data points in the training set in a similar way. Once we have computed the distances, we can find the nearest neighbor by selecting the data point with the smallest distance:

$$\hat{y} = y_i, \text{ where } i = \operatorname{argmin}_{j=1}^n d(x, X_j)$$

In this case, the nearest neighbor is the first data point, since it has the smallest distance to x . Therefore, the predicted class for the new data point x is $\hat{y} = y_1 = 0$.

This is transduction : No model was formed

B. KNN

KNN - K - nearest Neighbors - is a non parametric, instance based supervised learning algorithm. It is used for both classification and regression tasks.

The basic idea of KNN is to use the information from the K nearest Neighbors of a given data point

and make a prediction, given the current set of data that we work with.

Definition II.2 Mathematical definition of KNN is defined :

$$\hat{y} = \frac{1}{k} \sum_{i=1}^k y_i$$

For example, consider a dataset with two labels, 0 and 1, and a feature set containing two features, x_1 and x_2 . Let's say we want to predict the label of a new data

III. CONFORMAL PREDICTION

Definition III.1 Conformal prediction is a method of making predictions that is designed to provide a measure of confidence or uncertainty in the prediction. It is based on the idea of creating a prediction region, or "conformal prediction set," around a new data point, which contains a range of possible values for the prediction. The prediction is considered "valid" if the true value of the data point falls within this region.

The mathematical formula for conformal prediction is as follows:

$$P(\text{true label} \in \text{prediction set}) \geq 1 - \alpha$$

Here, $P(\text{true label} \in \text{prediction set})$ is the probability that the true label of the data point falls within the prediction set, and α is a user-specified confidence level, typically set to a small value such as 0.1 or 0.05.

Conformal prediction is used to provide a measure of uncertainty in predictions, which can be useful in situations where it is important to know how confident we are in a prediction. For example, in medical diagnosis, it may be important to know how certain we are that a patient has a particular disease, in order to make informed treatment decisions. In such cases, conformal prediction can be used to provide a range of possible diagnoses, along with a corresponding confidence level.

Conformal prediction is also used in situations where the data is noisy or uncertain, as it can provide a way to account for this uncertainty in the predictions. This is especially useful when working with complex or high-dimensional data, where the relationship between the features and the target variable may be difficult to model accurately.

a) Properties:

- Guaranteed validity - with respect to probability of error
- Commonplace in stats - Confidence Intervals, prediction intervals.

A. Explanation of Algo

The high level idea behind conformal prediction is to create a prediction set or region, around a new data point that contains a range of possible values for the prediction. This prediction set is constructed using the conformity scores and the p_value of the data points of the training and test values

- Given a training set of data points z_1, z_2, \dots, z_n and a new test sample x^* , we first compute the conformity scores, or p-values, for each possible label y for the test sample.
- The conformity score for a particular label y is calculated as the number of training data points with conformity scores less than or equal to the conformity score of the test sample, divided by the total number of training data points plus one. This is represented by the formula:

$$p(y) = \frac{i = 1, \dots, n + 1 | \alpha_y^i \leq \alpha_y^{n+1}}{n + 1}$$

where $\alpha_y^1, \dots, \alpha_y^n, \alpha_y^{n+1}$ are the conformity scores corresponding to $z_1, \dots, z_n, (x^*, y)$.

- Given a significance level $\epsilon > 0$ (our target probability of error), we can then compute the corresponding prediction set $\Gamma_\epsilon = \{y \in Y | p(y) > \epsilon\}$. This is the set of labels that are considered "valid" predictions for the test sample, as they have a probability of error less than or equal to ϵ .
- To make a prediction for the test sample, we can choose the label with the highest conformity score, or we can select the label with the highest conformity score that falls within the prediction set Γ_ϵ . This allows us to balance the trade-off between accuracy and confidence in the prediction.

Overall, the goal of conformal prediction is to provide a measure of confidence or uncertainty in the prediction, by constructing a prediction set that contains a range of possible values for the

prediction, rather than a single point estimate. This can be useful in situations where it is important to know how confident we are in a prediction, or where the data is noisy or uncertain.

1) *Special Cases*: Most cases it is important to understand how strange our data is. This refers to how different a data point is from the rest of the training data. A data point with a high degree of strangeness is one that is significantly different from the other data points in the training set.

The special cases are of the following: Because the slides are shit

- If the test sample has the highest conformity score. Of all the data points in the training set. This means that it is the most strange data point. In this case the conformity score of the test sample will be $\frac{1}{n+1}$ Which is the smallest possible value.
- If the test sample has the second highest conformity score of all the data points in the training set, this means that it is the second most strange data. In this case it will have a conformity score of the test $\frac{2}{n+1}$
- If the test sample has the lowest conformity score of all the data points in the training set, this means that it is the most "conforming" data point. In this case, the conformity score of the test sample will be 1.

These special cases illustrate how the conformity score of the test sample is influenced by its relationship to the conformity scores of the training data points. In general, the higher the conformity score of the test sample, the more "strange" it is relative to the training data, and the lower the conformity score, the more "conforming" it is. This is important to consider when constructing the prediction set, as the conformity scores of the test sample and the training data points are used to determine which labels are considered "valid" predictions for the test sample.

2) *Assumptions towards Conformity*: The assumptions of machine learning refer to the assumptions that are made about the data and the learning process when building and using machine learning models. These assumptions can vary depending on the type of model and the specific problem being solved, but they generally include assumptions about the structure of the data, the relationship between the features and the target

variable, and the nature of the learning process itself.

Conformal prediction is a method of making predictions that is designed to provide a measure of confidence or uncertainty in the prediction. It is based on the idea of creating a prediction region, or "conformal prediction set," around a new data point, which contains a range of possible values for the prediction. The prediction is considered "valid" if the true value of the data point falls within this region.

Conformal prediction based on nearest neighbors is a specific approach to conformal prediction that uses the k-nearest neighbors (KNN) algorithm to find the nearest neighbors of a new data point and use their labels to make a prediction. The conformity scores of the data points are calculated using the KNN distances, and the prediction set is constructed based on these conformity scores.

The validity and efficiency of conformal predictors refer to the properties of the prediction set that is constructed around a new data point. Validity refers to the probability that the true label of the data point falls within the prediction set, which is typically set to a user-specified value such as 0.1 or 0.05. Efficiency refers to the size of the prediction set, with smaller prediction sets being considered more efficient.

The validity of conformal predictors is ensured by the property that the predictor makes a mistake with probability at most ϵ , provided that the labeled samples are independently and identically distributed (IID). This means that the probability of the true label falling outside of the prediction set is at most ϵ .

The efficiency of conformal predictors can be traded off against their validity. In general, smaller prediction sets are more efficient, but they may also be less accurate, as they may not contain a wide enough range of possible values for the prediction. On the other hand, larger prediction sets may be more accurate, but they may also be less efficient. It is important to consider both validity and efficiency when using conformal prediction, in order to find the optimal balance between these two properties.

B. Conformity Measure for KNN ($K=1$)

There are several different conformity measures that can be used with the 1-nearest neighbor (1-NN) algorithm for conformal prediction. Here is an explanation of two of these measures:

The distance to the nearest sample of a different class: This conformity measure is based on the idea that a data point that is more similar to data points from a different class is more likely to be classified correctly. In other words, if the nearest neighbor of a data point belongs to a different class, this may be a good indication that the data point is correctly classified.

The conformity score for this measure can be calculated as the distance between the data point and its nearest neighbor of a different class. For example, if the data point is classified as class A and its nearest neighbor is classified as class B, the conformity score would be the distance between the two points.

This conformity measure can be represented mathematically as follows:

$$\alpha = d(\mathbf{x}, \mathbf{x}_{\text{NN}}^{\text{diff}})$$

where α is the conformity score, \mathbf{x} is the data point, and $\mathbf{x}_{\text{NN}}^{\text{diff}}$ is the nearest neighbor of a different class.

One over the distance to the nearest sample of the same class: This conformity measure is based on the opposite idea, that a data point that is more similar to data points from the same class is more likely to be classified correctly. In other words, if the nearest neighbor of a data point belongs to the same class, this may be a good indication that the data point is correctly classified.

The conformity score for this measure can be calculated as the reciprocal of the distance between the data point and its nearest neighbor of the same class. For example, if the data point is classified as class A and its nearest neighbor is also classified as class A, the conformity score would be the reciprocal of the distance between the two points.

This conformity measure can be represented mathematically as follows:

$$\alpha = \frac{1}{d(\mathbf{x}, \mathbf{x}_{\text{NN}}^{\text{same}})}$$

where α is the conformity score, \mathbf{x} is the data point, and $\mathbf{x}_{\text{NN}}^{\text{same}}$ is the nearest neighbor of the same class.

Both of these conformity measures are based on the distance between the data point and its nearest neighbors, with the first measure using the distance to the nearest neighbor of a different class and the second measure using the distance to the nearest neighbor of the same class.

The first conformity measure, which uses the distance to the nearest neighbor of a different class, may be more suitable in situations where the classes are well-separated and the data points from different classes are significantly different from one another. In such cases, the distance to the nearest neighbor of a different class may provide a good indication of the confidence of the prediction.

On the other hand, the second conformity measure, which uses the distance to the nearest neighbor of the same class, may be more suitable in situations where the classes are more overlapping and the data points from different classes are more similar to one another. In such cases, the distance to the nearest neighbor of the same class may provide a better indication of the confidence of the prediction.

C. Ranking Algo

It is a stupid algo, as simple as it May look, for some reason it can become quite jaring to figure out. This should debunk majority of its bullshit

Definition III.2 A ranking algorithm is a method for ordering a list of items according to some criterion. The criterion could be based on a variety of factors, such as relevance, popularity, or quality. We will use this for the p_value, to rank items that will have the same probability

1) *Example:* Imagine we have a list of students and their grades on a test, and we want to rank the students by their grades. The list looks like this:

Student	Grade
Alice	90
Bob	80
Charlie	90
Dave	70
Eve	80
Frank	60

To deal with duplicate values, we can first sort the list in descending order of grades. This will group the students with the same grade together:

Student	Grade
Alice	90
Charlie	90
Bob	80
Eve	80
Dave	70
Frank	60

Rank	Student	Grade
1	Alice	90
2	Charlie	90
3	Bob	80
3	Eve	80
4	Dave	70
5	Frank	60

- Sort the list in descending order of grades.
- Initialize the rank to 1.
- For each student in the list:
 - If the student has a higher grade than the previous student, increment the rank by 1.
 - If the student has the same grade as the previous student, do not increment the rank.
 - Assign the rank to the student.
- return Ranked List

D. Example

Take our original data Where we want to calculate the following

- P_value
- Conformity Score

Sample	Label
(0,3)	+1
(2,2)	+1
(3,3)	+1
(-1,1)	-1
(-1,-1)	-1
(0,1)	-1

For example, consider a dataset with two labels, 0 and 1, and a feature set containing two features, x_1 and x_2 . Where $x_1 = 0$ and $x_2 = 0$.

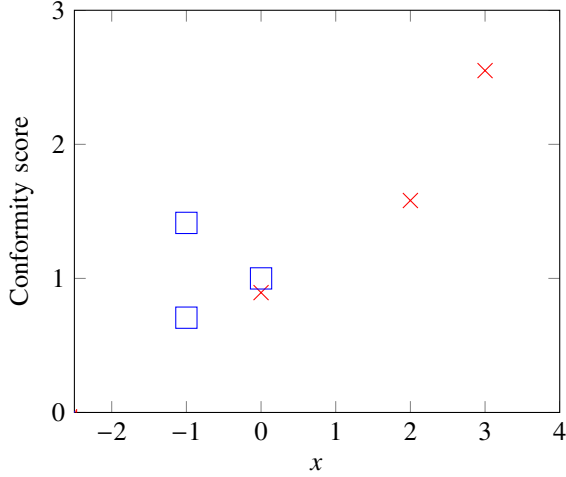
a) Steps: To calculate the conformity scores and p-values for the conformity measure "the distance to the nearest sample of a different class divided by the distance to the nearest sample of the same class," you can follow these steps:

- For each sample in the training set, calculate the conformity score for the given conformity measure. To do this, you will need to compute the distance to the nearest sample of a different class and the distance to the nearest sample of the same class. For example, for the positive sample (0,3), the conformity score would be calculated as $\frac{2}{\sqrt{4}+1} \approx 0.894$, where 2 is the distance to the nearest negative sample (-1,1) and $\sqrt{4}+1$ is the distance to the nearest positive sample (2,2).
- Repeat the calculation for the test sample, using the label you are interested in. For example, if you are interested in the p-value for the positive label, you would calculate the conformity score for the test sample (0,0) with the label +1. The conformity score for the test sample with the positive label would be calculated as $\frac{1}{\sqrt{4}+4} \approx 0.354$, where 1 is the distance to the nearest negative sample (-1,-1) and $\sqrt{4}+4$ is the distance to the nearest positive sample (3,3).
- Sort the conformity scores of the training samples, including the conformity score of the test sample, in ascending order.
- Count the number of conformity scores that are less than or equal to the conformity score of the test sample.
- Divide the number of conformity scores that are less than or equal to the conformity score of the test sample by the total number of conformity scores plus one, to get the p-value.

b) Calculation:

Sample	Label	Conformity score
(0, 3)	+1	$\frac{2}{\sqrt{4+1}} \approx 0.894$
(2, 2)	+1	$\frac{\sqrt{4+1}}{\sqrt{1+1}} \approx 1.581$
(3, 3)	+1	$\frac{\sqrt{9+4}}{\sqrt{1+1}} \approx 2.550$
(-1, 1)	-1	$\frac{\sqrt{1+1}}{1} \approx 1.414$
(-1, -1)	-1	$\frac{\sqrt{1+1}}{2} \approx 0.707$
(0, 1)	-1	$\frac{1}{1} = 1$
(0, 0)	+1	(?) $\frac{1}{\sqrt{4+4}} \approx 0.354$

Graphical Representation :



Sample	Label	Conformity score difference
(0, 3)	+1	0.894
(2, 2)	+1	1.581
(3, 3)	+1	2.550
(1, 1)	1	1.414
(1, 1)	1	0.707
(0, 1)	1	1

The conformity score of a sample can also be represented using the following equation:

$$\alpha = d(\mathbf{x}, \mathbf{x}_{\text{NN}}^{\text{diff}})$$

In this equation, α represents the conformity score of the sample, \mathbf{x} represents the coordinates of the sample, and $d(\mathbf{x}, \mathbf{x}_{\text{NN}}^{\text{diff}})$ represents the distance between the sample and its nearest neighbor of a different class. The nearest neighbor of a different class is represented by $\mathbf{x}_{\text{NN}}^{\text{diff}}$.

For example, to calculate the conformity score of the first positive sample (0, 3), you can use the following equation:

$$\alpha = d((0, 3), (-1, -1))$$

This equation calculates the distance between the first positive sample (0, 3) and its nearest neighbor of a different class (-1, -1). The distance can be calculated using the Euclidean distance formula:

$$d((0, 3), (-1, -1)) = \sqrt{(-1 - 0)^2 + (-1 - 3)^2} = 2$$

Substituting this value into the conformity score equation gives:

$$\alpha = \frac{2}{\sqrt{4+1}} = 0.894$$

This calculation can be repeated for each sample in the dataset to obtain the conformity scores for all samples.

c) *Prediction and how does ranking work* ? : To calculate the p-value using the ranking method, you need to rank all conformity scores in ascending order. The p-value for a particular label is then calculated by dividing the rank of the conformity score for that label by the total number of samples in the dataset, including the test sample.

For example, to calculate the p-value for the label +1 in the given dataset, you can first rank all conformity scores in ascending order as follows: 0.707, 0.894, 1.414, 1.581, 2.550, 0.354

The conformity score for the label +1 (0.354) is the sixth smallest conformity score. Therefore, the p-value for the label +1 is calculated as follows:

$$p(+1) = \frac{6}{7} = 0.857$$

Similarly, the p-value for the label -1 is calculated as follows:

$$p(-1) = \frac{1}{7} = 0.143$$

The p-values for all labels can be calculated in this way, and the prediction set can then be determined by comparing the p-values to the significance level ϵ . If a p-value is greater than ϵ , the corresponding label is included in the prediction set. For example, if the significance level is set to $\epsilon = 0.1$, then the prediction set for the given dataset would be $\Gamma_{0.1} = -1$, since the p-value for the label -1 (0.143) is greater than ϵ .

d) *Ranking*: In the ranking method, the conformity score for the test sample is compared to the conformity scores for all other samples in the dataset. If the conformity score for the test sample is the smallest (or largest), it gets a rank of 1 (or the total number of samples). If the conformity score for the test sample is the second smallest (or second largest), it gets a rank of 2 (or the total number of samples minus 1). This process continues until the conformity score for the test sample is ranked.

In the given dataset, the conformity score for the test sample (0.354) is the sixth smallest conformity score. Therefore, it gets a rank of 6. The total number of samples in the dataset, including the test sample, is 7. Therefore, the p-value for the label +1 is calculated as follows:

$$p(+1) = \frac{6}{7} = 0.857$$

Similarly, the conformity score for the label -1 (0.707) is the smallest conformity score and gets a rank of 1. The p-value for the label -1 is calculated as follows:

$$p(-1) = \frac{1}{7} = 0.143$$

This is why the p-values for the labels +1 and -1 are 6/7 and 1/7, respectively.