

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Jan 20, 2021

1 / 42

Logistics

Outline

- 1 Logistics
- 2 Recap
- 3 Classification and Nearest Neighbor Classifier (NNC)
- 4 Some theory on NNC

3 / 42

Outline

- 1 Logistics
- 2 Recap
- 3 Classification and Nearest Neighbor Classifier (NNC)
- 4 Some theory on NNC

2 / 42

Logistics

Homeworks

- HW 0 is due on Friday. It will not be graded, only to get everyone familiar with the submission mechanism.
- HW 1 will be released on Friday (01/22/2021). Starting HW 1 assignments will be graded.

4 / 42

Outline

- 1 Logistics
- 2 **Recap**
- 3 Classification and Nearest Neighbor Classifier (NNC)
- 4 Some theory on NNC

5 / 42

Last Class: Foundations of ML

- We discussed different flavors of learning problems
- Tools from probability, information theory, and optimization
- **Today:** We will start our journey of Supervised learning starting with classification.

6 / 42

Outline

- 1 Logistics
- 2 Recap
- 3 **Classification and Nearest Neighbor Classifier (NNC)**
 - Intuitive example
 - General setup for classification
 - Algorithm
 - How to measure performance
 - Variants, Parameters, and Tuning
 - Summary
- 4 Some theory on NNC

7 / 42

Recognizing flowers

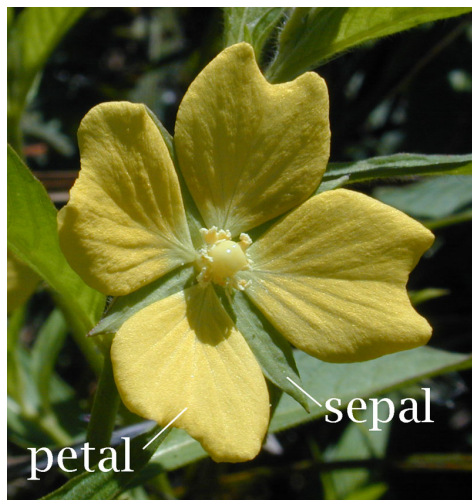
Types of Iris: setosa, versicolor, and virginica



8 / 42

Measuring the properties of the flowers

Features and attributes: the widths and lengths of sepal and petal



9 / 42

Often, data is conveniently organized as a table

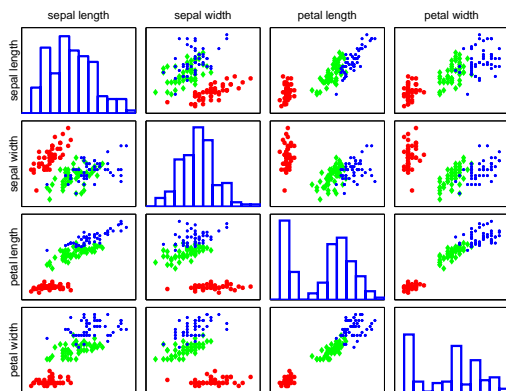
Fisher's Iris Data				
Sepal length ↕	Sepal width ↕	Petal length ↕	Petal width ↕	Species ↕
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>
4.4	2.9	1.4	0.2	<i>I. setosa</i>
4.9	3.1	1.5	0.1	<i>I. setosa</i>

10 / 42

Pairwise scatter plots of 131 flower specimens

Visualization of data helps identify the right learning model to use

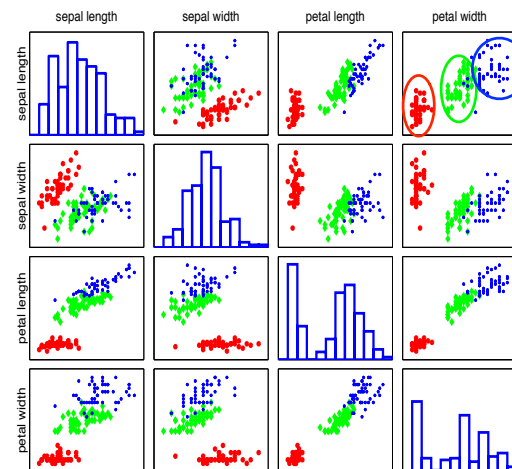
Each colored point is a flower specimen: *setosa* (red), *versicolor* (green), *virginica* (blue)



11 / 42

Different types seem well-clustered and separable

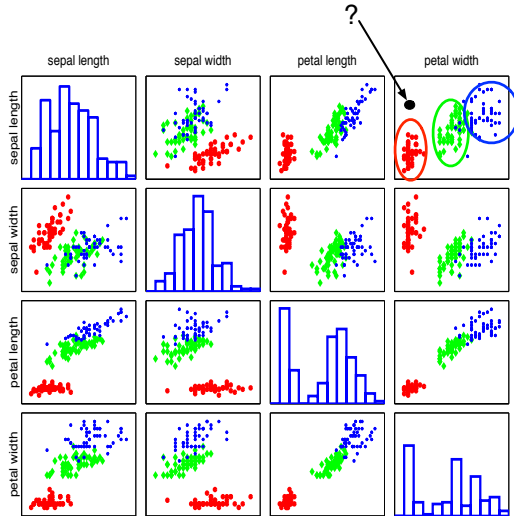
Using two features: petal width and sepal length



12 / 42

Labeling an unknown flower type

Closer to red cluster: so labeling it as **setosa**



13 / 42

General setup for multi-class classification

Training data (set)

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$
- Each $x_n \in \mathbb{R}^D$ is called a feature vector.
- Each $y_n \in [C] = \{1, 2, \dots, C\}$ is called a label/class/category.
- They are used to learn a *classifier* $f : \mathbb{R}^D \rightarrow [C]$ for future prediction.

Special case: binary classification

- Number of classes: $C = 2$
- Conventional labels: $\{0, 1\}$ or $\{-1, +1\}$

14 / 42

Nearest neighbor classification (NNC)

The index of the **nearest neighbor** of a point x is

$$\text{nn}(x) = \underset{n \in [N]}{\operatorname{argmin}} \|x - x_n\|_2 = \underset{n \in [N]}{\operatorname{argmin}} \sqrt{\sum_{d=1}^D (x_d - x_{nd})^2}$$

where $\|\cdot\|_2$ is the ℓ_2 /Euclidean distance.

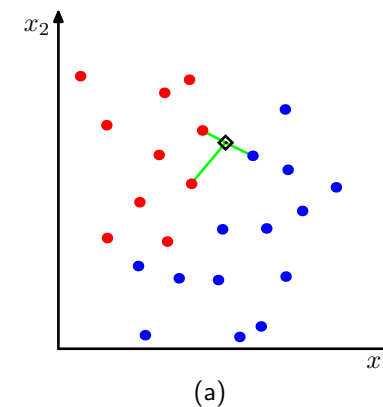
Classification rule

$$f(x) = y_{\text{nn}(x)}$$

15 / 42

Visual example

In this 2-dimensional example, the nearest point to x is a **red training instance**, thus, x will be labeled as **red**.



16 / 42

Example: classify Iris with two features

Training data

ID (n)	petal width (x_1)	sepal length (x_2)	category (y)
1	0.2	5.1	setoas
2	1.4	7.0	versicolor
3	2.5	6.7	virginica
\vdots	\vdots	\vdots	

Flower with unknown category

petal width = 1.8 and sepal length = 6.4 (i.e. $\mathbf{x} = (1.8, 6.4)$)

Calculating distance $\|\mathbf{x} - \mathbf{x}_n\|_2 = \sqrt{(x_1 - x_{n1})^2 + (x_2 - x_{n2})^2}$

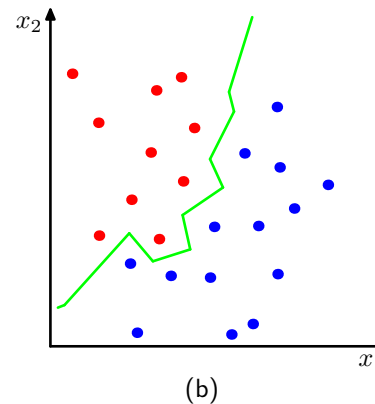
ID	distance
1	1.75
2	0.72
3	0.76

Thus, the category is *versicolor*.

17 / 42

Decision boundary

For every point in the space, we can determine its label using the NNC rule. This gives rise to a *decision boundary* that partitions the space into different regions.



18 / 42

Is NNC doing the right thing for us?

Intuition

We should compute **accuracy** (A) — the percentage of data points being correctly classified, or the **error rate** (ϵ) — the percentage of data points being incorrectly classified. (accuracy + error rate = 1)

Defined on the training data set

$$A^{\text{TRAIN}} = \frac{1}{N} \sum_n \mathbb{I}[f(\mathbf{x}_n) == y_n], \quad \epsilon^{\text{TRAIN}} = \frac{1}{N} \sum_n \mathbb{I}[f(\mathbf{x}_n) \neq y_n]$$

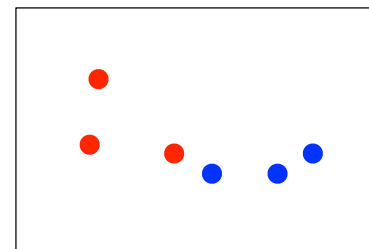
where $\mathbb{I}[\cdot]$ is the indicator function.

Is this the right measure?

19 / 42

Example

Training data



What are A^{TRAIN} and ϵ^{TRAIN} ?

$$A^{\text{TRAIN}} = 100\%, \quad \epsilon^{\text{TRAIN}} = 0\%$$

For every training data point, its nearest neighbor is itself.

20 / 42

Test Error

Does it mean nearest neighbor is a very good algorithm?

Not really, having zero training error is simple!

We should care about accuracy when predicting unseen data

Test/Evaluation data

- $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- A fresh dataset, *not* overlap with training set.
- Test accuracy and test error

$$A^{\text{TEST}} = \frac{1}{M} \sum_m \mathbb{I}[f(\mathbf{x}_m) == y_m], \quad \epsilon^{\text{TEST}} = \frac{1}{M} \sum_m \mathbb{I}[f(\mathbf{x}_m) \neq y_m]$$

- Good measurement of a classifier's performance

Variant 1: measure nearness with other distances

Previously, we use the Euclidean distance

$$\text{nn}(\mathbf{x}) = \underset{n \in [\mathbf{N}]}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_n\|_2$$

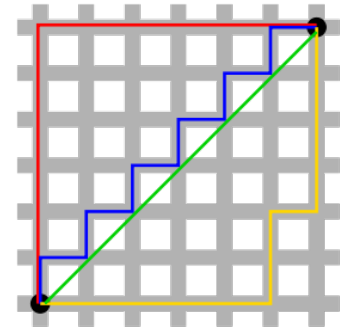
Many other alternative distances

E.g., the following L_1 distance (i.e., city block distance, or Manhattan distance)

$$\|\mathbf{x} - \mathbf{x}_n\|_1 = \sum_{d=1}^D |x_d - x_{nd}|$$

More generally, L_p distance (for $p \geq 1$):

$$\|\mathbf{x} - \mathbf{x}_n\|_p = \left(\sum_d |x_d - x_{nd}|^p \right)^{1/p}$$



Green line is Euclidean distance.
Red, Blue, and Yellow lines are L_1 distance

Variant 2: K-nearest neighbor (KNN)

Increase the number of nearest neighbors to use?

- 1-nearest neighbor: $\text{nn}_1(\mathbf{x}) = \underset{n \in [\mathbf{N}]}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_n\|_2$
- 2-nearest neighbor: $\text{nn}_2(\mathbf{x}) = \underset{n \in [\mathbf{N}] \setminus \{\text{nn}_1(\mathbf{x})\}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_n\|_2$
- 3-nearest neighbor: $\text{nn}_3(\mathbf{x}) = \underset{n \in [\mathbf{N}] \setminus \{\text{nn}_1(\mathbf{x}), \text{nn}_2(\mathbf{x})\}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_n\|_2$

The set of K-nearest neighbor

$$\text{knn}(\mathbf{x}) = \{\text{nn}_1(\mathbf{x}), \text{nn}_2(\mathbf{x}), \dots, \text{nn}_K(\mathbf{x})\}$$

Note: we have

$$\|\mathbf{x} - \mathbf{x}_{\text{nn}_1(\mathbf{x})}\|_2 \leq \|\mathbf{x} - \mathbf{x}_{\text{nn}_2(\mathbf{x})}\|_2 \leq \dots \leq \|\mathbf{x} - \mathbf{x}_{\text{nn}_K(\mathbf{x})}\|_2$$

How to classify with K neighbors?

Classification rule

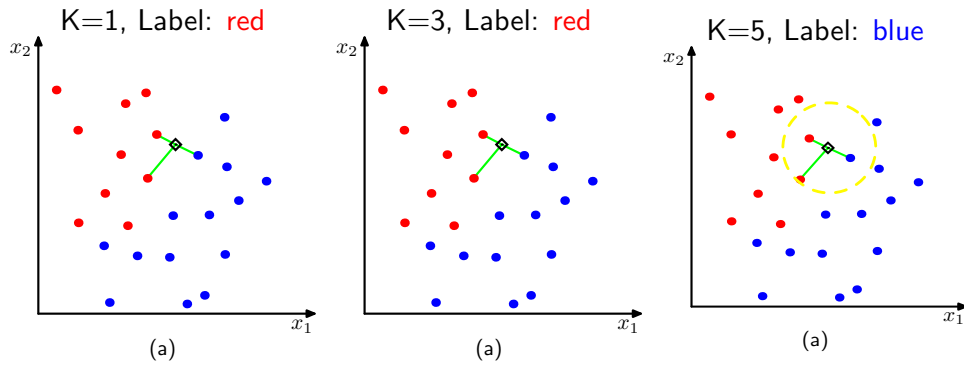
- Every neighbor votes: naturally \mathbf{x}_n votes for its label y_n .
- Aggregate everyone's vote on a class label c

$$v_c = \sum_{n \in \text{knn}(\mathbf{x})} \mathbb{I}(y_n == c), \quad \forall c \in [\mathbf{C}]$$

- Predict with the majority

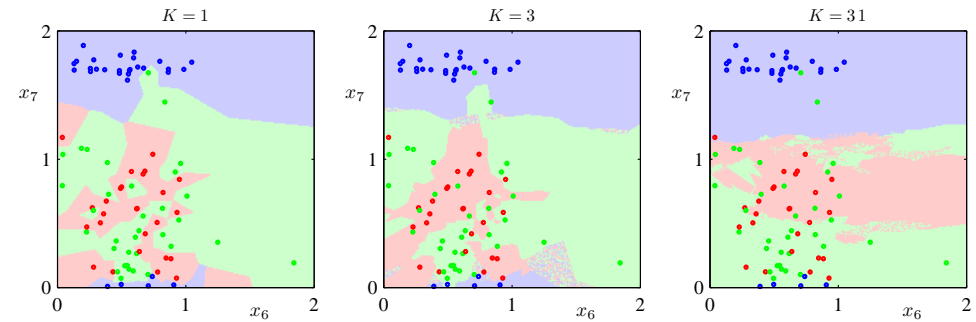
$$f(\mathbf{x}) = \underset{c \in [\mathbf{C}]}{\operatorname{argmax}} v_c$$

Example



25 / 42

Decision boundary



When K increases, the decision boundary becomes smoother.

What happens when $K = N$?

26 / 42

Variant 3: Preprocessing data

One issue of NNC: *distances depend on units of the features!*

One solution: preprocess data so it looks more “normalized”.

Example:

- compute the means and standard deviations in each feature

$$\bar{x}_d = \frac{1}{N} \sum_n x_{nd}, \quad s_d^2 = \frac{1}{N-1} \sum_n (x_{nd} - \bar{x}_d)^2$$

- Scale the feature accordingly

$$x_{nd} \leftarrow \frac{x_{nd} - \bar{x}_d}{s_d}$$

Many other ways of normalizing data.

27 / 42

Which variants should we use?

Hyper-parameters in NNC

- The distance measure (e.g. the parameter p for L_p norm)
- K (i.e. how many nearest neighbor?)
- Different ways of preprocessing

Most algorithms have hyper-parameters. Tuning them is a significant part of applying an algorithm.

28 / 42

Tuning via a development dataset

Training data

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- They are used to learn $f(\cdot)$

Test data

- M samples/instances: $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used to evaluate how well $f(\cdot)$ will do.

Development/Validation data

- L samples/instances: $\mathcal{D}^{\text{DEV}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_L, y_L)\}$
- They are used to optimize hyper-parameter(s).

These three sets should *not* overlap!

29 / 42

Recipe

- For each possible value of the hyperparameter (e.g. $K = 1, 3, \dots$)
 - Train a model using $\mathcal{D}^{\text{TRAIN}}$
 - Evaluate the performance of the model on \mathcal{D}^{DEV}
- Choose the model with the best performance on \mathcal{D}^{DEV}
- Evaluate the model on $\mathcal{D}^{\text{TEST}}$

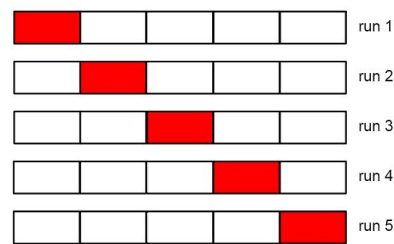
30 / 42

S-fold Cross-validation

What if we do not have a development set?

- Split the training data into S equal parts.
- Use each part *in turn* as a development dataset and use the others as a training dataset.
- Choose the hyper-parameter leading to best *average* performance.

S = 5: 5-fold cross validation



Special case: $S = N$, called leave-one-out.

31 / 42

Cross-validation recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$.
- For each possible value of the hyper-parameter (e.g. $K = 1, 3, \dots$)
 - For every $s \in [S]$
 - Train a model using $\mathcal{D}_{\setminus s}^{\text{TRAIN}} = \mathcal{D}^{\text{TRAIN}} - \mathcal{D}_s^{\text{TRAIN}}$
 - Evaluate the performance of the model on $\mathcal{D}_s^{\text{TRAIN}}$
 - Average the S performance metrics
- Choose the hyper-parameter with the best averaged performance
- **Use the best hyper-parameter to train a model using all $\mathcal{D}^{\text{train}}$**
- Evaluate the model on $\mathcal{D}^{\text{TEST}}$

32 / 42

Summary

Advantages of NNC

- Simple, easy to implement (widely used in practice)

Disadvantages of NNC

- Computationally intensive for large-scale problems: $O(ND)$ for each prediction *naively*. Here, N is the cardinality of the training set and D is the dimension of the training example.
- Need to “*carry*” the training data around. This type of method is called *nonparametric*.
- Choosing the right hyper-parameters can be involved.

33 / 42

Outline

- 1 Logistics
- 2 Recap
- 3 Classification and Nearest Neighbor Classifier (NNC)
- 4 Some theory on NNC
 - Step 1: Expected risk
 - Step 2: The ideal classifier
 - Step 3: Comparing NNC to the ideal classifier

35 / 42

Summary

Typical steps of developing a machine learning system:

- Collect data, split into training, development, and test sets.
- Train a model with a machine learning algorithm. Most often we apply cross-validation to tune hyper-parameters.
- Evaluate using the test data and report performance.
- Use the model to predict future/make decisions.

34 / 42

How good is NNC really?

To answer this question, we proceed in 3 steps

- 1 Define *more carefully* a performance metric for a classifier.
- 2 Hypothesize an ideal classifier - *the best possible one*.
- 3 Compare NNC to the ideal one.

36 / 42

Why does test error make sense?

Test error makes sense only when training set and test set are correlated.

Most standard assumption: every data point (\mathbf{x}, y) (from $\mathcal{D}^{\text{TRAIN}}$, \mathcal{D}^{DEV} , or $\mathcal{D}^{\text{TEST}}$) is an *independent and identically distributed (i.i.d.)* sample of an unknown joint distribution \mathcal{P} .

- often written as $(\mathbf{x}, y) \stackrel{i.i.d.}{\sim} \mathcal{P}$

Test error of a fixed classifier is therefore a *random variable*.

Need a more “certain” measure of performance (so it’s easy to compare different classifiers for example).

37 / 42

Expected risk

More generally, for a loss function $L(y', y)$,

- e.g. $L(y', y) = \mathbb{I}[y' \neq y]$, called *0-1 loss*. **Default**
- many more other losses as we will see.

the *expected risk* of f is defined as

$$R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} L(f(\mathbf{x}), y).$$

For *0-1 loss* we have

$$R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \mathbb{I}[y' \neq y]$$

39 / 42

Expected error

What about the **expectation** of this random variable?

$$\mathbb{E}[\epsilon^{\text{TEST}}] = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{(\mathbf{x}_m, y_m) \sim \mathcal{P}} \mathbb{I}[f(\mathbf{x}_m) \neq y_m] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} \mathbb{I}[f(\mathbf{x}) \neq y]$$

- i.e. the expected error/mistake of f

Test error is a proxy of expected error. *The larger the test set, the better the approximation.*

What about the expectation of training error? Is training error a good proxy of expected error?

38 / 42

Bayes optimal classifier

What should we predict for \mathbf{x} , *knowing* $\mathcal{P}(y|\mathbf{x})$?

Bayes optimal classifier: $f^*(\mathbf{x}) = \operatorname{argmax}_{c \in [C]} \mathcal{P}(c|\mathbf{x})$.

The optimal risk: $R(f^*) = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_\mathbf{x}} [1 - \max_{c \in [C]} \mathcal{P}(c|\mathbf{x})]$ where $\mathcal{P}_\mathbf{x}$ is the marginal distribution of \mathbf{x} .

That is we have $R(f^*) \leq R(f)$ for any f . **Verify!**

For special case $C = 2$, let $\eta(\mathbf{x}) = \mathcal{P}(0|\mathbf{x})$, then

$$\begin{aligned} R(f^*) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_\mathbf{x}} [\mathbb{E}_{\mathbf{y}|\mathbf{x}} [\mathbb{I}_{f^*(\mathbf{x}) \neq y}]] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_\mathbf{x}} [\eta(\mathbf{x}) \mathbb{I}_{f^*(\mathbf{x})=1} + (1 - \eta(\mathbf{x})) \mathbb{I}_{f^*(\mathbf{x})=0}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{P}_\mathbf{x}} [\min\{\eta(\mathbf{x}), 1 - \eta(\mathbf{x})\}], \end{aligned}$$

40 / 42

Comparing NNC to Bayes optimal classifier

Come back to the question: how good is NNC?

Theorem (Cover and Hart, 1967)

Let f_N be the 1-nearest neighbor binary classifier using N training data points, we have (under mild conditions)

$$R(f^*) \leq \lim_{N \rightarrow \infty} \mathbb{E}[R(f_N)] \leq 2R(f^*)$$

i.e., expected risk of NNC in the limit is at most twice of the best possible.

A pretty strong guarantee.

In particular, $R(f^*) = 0$ implies $\mathbb{E}[R(f_N)] \rightarrow 0$.

Proof sketch

Fact: $x_{nn(x)} \rightarrow x$ as $N \rightarrow \infty$ with probability 1

$$\begin{aligned} \mathbb{E}[R(f_N)] &= \mathbb{E}[\mathbb{E}_{(x,y) \sim \mathcal{P}} \mathbb{I}[f_N(x) \neq y]] \\ &\rightarrow \mathbb{E}_{x \sim \mathcal{P}_x} \mathbb{E}_{y, y' \stackrel{i.i.d.}{\sim} \mathcal{P}(\cdot|x)} [\mathbb{I}[y' \neq y]] \\ &= \mathbb{E}_{x \sim \mathcal{P}_x} \mathbb{E}_{y, y' \stackrel{i.i.d.}{\sim} \mathcal{P}(\cdot|x)} [\mathbb{I}[y' = 0 \text{ and } y = 1] + \mathbb{I}[y' = 1 \text{ and } y = 0]] \\ &= \mathbb{E}_{x \sim \mathcal{P}_x} [\eta(x)(1 - \eta(x)) + (1 - \eta(x))\eta(x)] \\ &= 2\mathbb{E}_{x \sim \mathcal{P}_x} [\eta(x)(1 - \eta(x))] \\ &\leq 2\mathbb{E}_{x \sim \mathcal{P}_x} [\min\{\eta(x), (1 - \eta(x))\}] \\ &= 2R(f^*) \end{aligned}$$