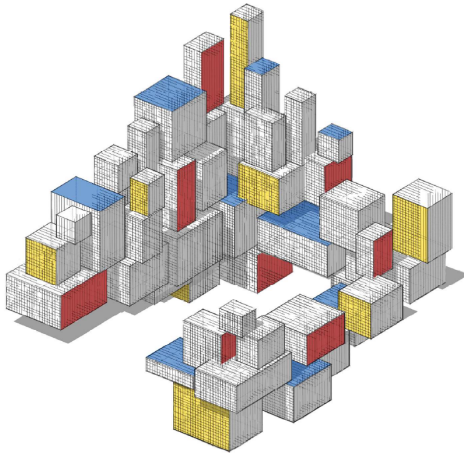


Logistic Regression and NN Classification



Purpose

In this lecture we will discuss:

- The multi-logit model
- Softmax classification
- Nearest Neighbor classification

Logistic Regression and Softmax Classification

We introduced the logistic (logit) regression model as a generalized linear model where, conditional on a p -dimensional feature vector \mathbf{x} , the random response $Y \sim \text{Ber}(h(\mathbf{x}^\top \boldsymbol{\beta}))$, with

$$h(u) = 1/(1 + e^{-u}).$$

The parameter $\boldsymbol{\beta}$ is learned from the training data by maximizing the likelihood of the training responses or, equivalently, by minimizing the **cross-entropy training loss**:

$$-\frac{1}{n} \sum_{i=1}^n \ln g(y_i | \boldsymbol{\beta}, \mathbf{x}_i),$$

where $g(y = 1 | \boldsymbol{\beta}, \mathbf{x}) = \frac{1}{1+e^{-\mathbf{x}^\top \boldsymbol{\beta}}}$ and $g(y = 0 | \boldsymbol{\beta}, \mathbf{x}) = \frac{e^{-\mathbf{x}^\top \boldsymbol{\beta}}}{1+e^{-\mathbf{x}^\top \boldsymbol{\beta}}}$.

In particular, we have

$$\ln \frac{g(y = 1 | \boldsymbol{\beta}, \mathbf{x})}{g(y = 0 | \boldsymbol{\beta}, \mathbf{x})} = \mathbf{x}^\top \boldsymbol{\beta}. \quad (1)$$

Logit Classification

Hence, the decision boundary $\{\mathbf{x} : g(y = 0 | \boldsymbol{\beta}, \mathbf{x}) = g(y = 1 | \boldsymbol{\beta}, \mathbf{x})\}$ is the hyperplane $\mathbf{x}^\top \boldsymbol{\beta} = 0$.

If the constant feature is considered separately, that is $\mathbf{x} = [1, \tilde{\mathbf{x}}^\top]^\top$, then the boundary is an affine hyperplane in $\tilde{\mathbf{x}}$.

Suppose that training on $\tau = \{(\mathbf{x}_i, y_i)\}$ yields the estimate $\hat{\boldsymbol{\beta}}$ with the corresponding learner $g_\tau(y = 1 | \mathbf{x}) = 1/(1 + e^{-\mathbf{x}^\top \hat{\boldsymbol{\beta}}})$.

The learner can be used as a **pre-classifier** from which we obtain the classifier $\mathbb{I}\{g_\tau(y = 1 | \mathbf{x}) > 1/2\}$ or, equivalently,

$$\hat{y} := \operatorname{argmax}_{j \in \{0,1\}} g_\tau(y = j | \mathbf{x}).$$

Multi-Logit Classification

The classification methodology for the logit model can be generalized to the **multi-logit** model where the response takes values in the set $\{0, \dots, c - 1\}$.

The key idea is to replace (1) with

$$\ln \frac{g(y = j \mid \mathbf{W}, \mathbf{b}, \mathbf{x})}{g(y = 0 \mid \mathbf{W}, \mathbf{b}, \mathbf{x})} = \mathbf{x}^\top \boldsymbol{\beta}_j, \quad j = 1, \dots, c - 1, \quad (2)$$

where the matrix $\mathbf{W} \in \mathbb{R}^{(c-1) \times (p-1)}$ and vector $\mathbf{b} \in \mathbb{R}^{c-1}$ reparameterize all $\boldsymbol{\beta}_j \in \mathbb{R}^p$ such that (recall $\mathbf{x} = [1, \tilde{\mathbf{x}}^\top]^\top$):

$$\mathbf{W} \tilde{\mathbf{x}} + \mathbf{b} = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_{c-1}]^\top \mathbf{x}.$$

Thus, the random response Y is assumed to have a conditional probability distribution for which the **log-odds ratio** with respect to class j and a “reference” class (in this case 0) is **linear**.

Multi-Logit Classification

The model (2) completely specifies the distribution of Y , namely:

$$g(y \mid \mathbf{W}, \mathbf{b}, \mathbf{x}) = \frac{\exp(z_{y+1})}{\sum_{k=1}^c \exp(z_k)}, \quad y = 0, \dots, c-1,$$

where z_1 is an arbitrary constant, say 0, corresponding to the “reference” class $y = 0$, and

$$[z_2, \dots, z_c]^\top := \mathbf{W} \tilde{\mathbf{x}} + \mathbf{b}.$$

Note: $g(y \mid \mathbf{W}, \mathbf{b}, \mathbf{x})$ is the $(y+1)$ -st component of $\mathbf{a} = \text{softmax}(\mathbf{z})$, where

$$\text{softmax} : \mathbf{z} \mapsto \frac{\exp(z)}{\sum_k \exp(z_k)}$$

is the **softmax** function and $\mathbf{z} = [z_1, \dots, z_c]^\top$.

Multi-Logit Classification

Finally, we can write the classifier as

$$\hat{y} = \operatorname{argmax}_{j \in \{0, \dots, c-1\}} a_{j+1}.$$

In summary, we have the sequence of mappings transforming the input \mathbf{x} into the output \hat{y} :

$$\mathbf{x} \rightarrow \mathbf{W}\tilde{\mathbf{x}} + \mathbf{b} \rightarrow \operatorname{softmax}(\mathbf{z}) \rightarrow \operatorname{argmax}_{j \in \{0, \dots, c-1\}} a_{j+1} \rightarrow \hat{y}.$$

In the context of neural networks, \mathbf{W} is called a **weight** matrix and \mathbf{b} is called a **bias** vector.

Training Loss

The parameters \mathbf{W} and \mathbf{b} have to be learned from the training data, which involves minimization of the (supervised version of) the **cross-entropy training loss**:

$$\frac{1}{n} \sum_{i=1}^n \text{Loss}(f(y_i | \mathbf{x}_i), g(y_i | \mathbf{W}, \mathbf{b}, \mathbf{x}_i)) = -\frac{1}{n} \sum_{i=1}^n \ln g(y_i | \mathbf{W}, \mathbf{b}, \mathbf{x}_i).$$

Using the softmax function, the **cross-entropy** loss can be simplified to:

$$\text{Loss}(f(y | \mathbf{x}), g(y | \mathbf{W}, \mathbf{b}, \mathbf{x})) = -z_{y+1} + \ln \sum_{k=1}^c \exp(z_k).$$

K -Nearest Neighbors Classification

Let $\tau = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ be the training set, with $y_i \in \{0, \dots, c-1\}$, and let \mathbf{x} be a new feature vector.

Define $\mathbf{x}_{(1)}, \mathbf{x}_{(2)}, \dots, \mathbf{x}_{(n)}$ as the feature vectors ordered by closeness to \mathbf{x} in some distance $\text{dist}(\mathbf{x}, \mathbf{x}_i)$, e.g., the Euclidean distance $\|\mathbf{x} - \mathbf{x}'\|$.

Let $\tau(\mathbf{x}) := \{(\mathbf{x}_{(1)}, y_{(1)}) \dots, (\mathbf{x}_{(K)}, y_{(K)})\}$ be the subset of τ that contains K feature vectors \mathbf{x}_i that are closest to \mathbf{x} .

Then the K -nearest neighbors classification rule classifies \mathbf{x} according to the most frequently occurring class labels in $\tau(\mathbf{x})$.

If two or more labels receive the same number of votes, the feature vector is classified by selecting one of these labels randomly with equal probability.

Voronoi Tessellation

For the case $K = 1$ the set $\tau(\mathbf{x})$ contains only one element, say (\mathbf{x}', y') , and \mathbf{x} is classified as y' . This divides the space into n regions

$$\mathcal{R}_i = \{\mathbf{x} : \text{dist}(\mathbf{x}, \mathbf{x}_i) \leq \text{dist}(\mathbf{x}, \mathbf{x}_j), j \neq i\}, \quad i = 1, \dots, n.$$

For a feature space \mathbb{R}^p with the Euclidean distance, this gives a **Voronoi tessellation** of the feature space.

The following figure and Python program show how nearest neighbor classification works for 80 randomly simulated points above and below the line $x_2 = x_1$. Points above the line $x_2 = x_1$ have label 0 and points below this line have label 1.

Example: Nearest Neighbor Classification

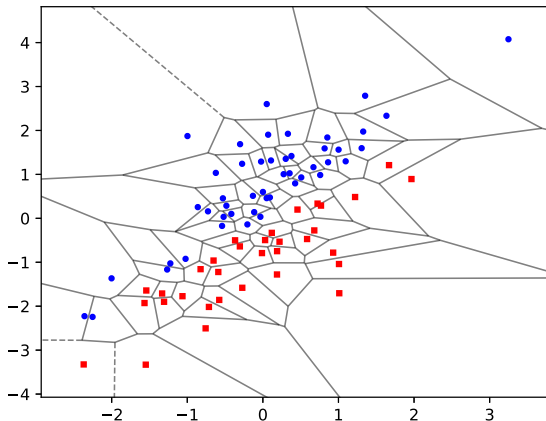


Figure: The 1-nearest neighbor algorithm divides up the space into Voronoi cells.

```
import numpy as np
from numpy.random import rand, randn
import matplotlib.pyplot as plt
from scipy.spatial import Voronoi, voronoi_plot_2d

np.random.seed(12345)
M = 80
x = randn(M,2)
y = np.zeros(M) # pre-allocate list

for i in range(M):
    if rand()<0.5:
        x[i,1], y[i] = x[i,0] + np.abs(randn()), 0
    else:
        x[i,1], y[i] = x[i,0] - np.abs(randn()), 1

vor = Voronoi(x)
plt_options = {'show_vertices':False, 'show_points':False,
               'line_alpha':0.5}
fig = voronoi_plot_2d(vor, **plt_options)
plt.plot(x[y==0,0], x[y==0,1], 'bo',
         x[y==1,0], x[y==1,1], 'rs', markersize=3)
```