

INT6151 Machine Learning

Lecture 3 - Classification

Ta Viet Cuong

VNU-UET

2024

Table of content

Multi-class logistics regression model

Multi-class classification

The optimal classifier

K nearest neighbor - KNN

Recap: Logistic Regression - Binary Classification

Data

$$x \in \mathbb{R}^d, y \in \{0, 1\}$$

Example: image $S \times S \rightarrow d = S^2 = 784$ (MNIST)

$$D = \{(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)\}$$

Model

$$f(x) = w^T x + w_0$$

$$Y|X = x \sim \text{Ber}(y|\sigma(f(x)))$$

Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Parameter

$$\theta = (w, w_0)$$

Recap: Logistic Regression - Binary cross entropy loss

Training

With MLE principle, calculate likelihood

$$L(w, w_0) = P(D) = \prod_{i=1}^n P(y_i | x_i) = \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i}$$

where $\mu_i = \sigma(f(x_i))$

Negative-loglikelihood (NLL)

$$\ell(w, w_0) = -\log L(w, w_0) = \sum_{i=1}^n -y_i \log \mu_i - (1 - y_i) \log(1 - \mu_i)$$

💡 Binary cross entropy loss function - BCE

Recap: Training Algorithm - Gradient Descent

TrainLR-GD(D, λ) $\rightarrow w, w_0$:

1. Initialize: $w = 0 \in \mathbb{R}^d, w_0 = 0$
2. Loop $epoch = 1, 2, \dots$
 - a. Calculate $\ell(w, w_0)$
 - b. Calculate derivative $\nabla_w \ell, \nabla_{w_0} \ell$
 - c. Update params

$$w \leftarrow w - \lambda \nabla_w \ell(w, w_0)$$

$$w_0 \leftarrow w_0 - \lambda \nabla_{w_0} \ell(w, w_0)$$

3. Stop when:
 - a. Epoch is large enough
 - b. The Loss function decrease negligible
 - c. The derivative is small enough $\|\nabla_w \ell\|, \nabla_{w_0} \ell$

Multi-class logistics regression model

Given a label set $\mathcal{Y} = \{1, 2, \dots, C\}$

Categorical distribution

A random variable $Y \sim \text{Cat}(y|\theta_1, \theta_2, \dots, \theta_C)$ means that

$$P(Y = c) = \theta_c, c = 1, 2, \dots, C$$

with θ_c is the probability of the category c and $\sum_c \theta_c = 1$

For example: A six-side dice with $C = 6, \theta_c = 1/6, \forall c$

$$P(Y = y) = \prod_{c=1}^C \theta_c^{\mathbb{I}(c=y)}$$

where $\mathbb{I}(c = y)$ is an indicator random function denoting whether $c = y$ or not.

A dataset example with categorical labels (1)

Iris dataset:

- ▶ Number of Instances: 150
- ▶ Number of Attributes: 4 (sepal length/width in centimeters, petal length/width in centimeters)
- ▶ Number of classes: 3

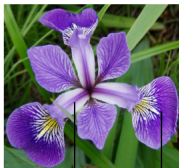
iris setosa



petal

sepal

iris versicolor



petal

sepal

iris virginica



petal

sepal

A dataset example with categorical labels (2)

MNIST dataset:

- ▶ Number of Instances: 60,000 training images and 10,000 testing images.
- ▶ Number of Attributes: 28×28
- ▶ Number of classes: 10



Figure: Sample images from MNIST test dataset (source: Wikipedia)

Multi-class logistic regression model

Model

Linear function with parameters $w_c \in \mathbb{R}^d$, $w_{c0} \in \mathbb{R}$, $c = 1, 2, \dots, C$

$$f_c(x) = w_c^T x + w_{c0}$$

with the application of Softmax function to convert a set of (z_1, z_2, \dots, z_C) to probabilities

$$\mathcal{S}(z_1, z_2, \dots, z_C) = \left[\frac{e^{z_c}}{\sum_{c'=1}^C e^{z_{c'}}} \right]_{c=1,2,\dots,C}$$

Probability model

$$Y|X = x \sim \text{Cat}(y|\mathcal{S}(f_1(x), f_2(x), \dots, f_C(x)))$$

where $f_c(x)$ is called as **logit**. With logistic regression, $f_c(x)$ are (simple) linear functions, more sophisticated methods make use of neural networks.

Multi-class logistic regression model

Inference

$h^{LR}(x)$: Choose c to maximize $h^{LR}(x)$,

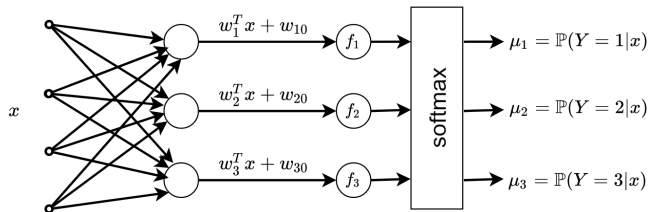


Figure: Multi-class logistic regression with Softmax

Multi-class logistic regression model

Training

The likelihood of the parameters with respect to the dataset:

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$\begin{aligned} L(\mathbf{W}) &= \mathbb{P}(D) = \prod_{i=1}^n \mathbb{P}(Y = y_i | x_i) \\ &= \prod_{i=1}^n \prod_{c=1}^C \mu_{ic}^{y_{ic}} \end{aligned}$$

where μ_{ic} is computed as:

$$\mu_{ic} = \frac{e^{f_c(x_i)}}{\sum_{c'=1}^C e^{f_{c'}(x_i)}}$$

and the one-hot encoding of the labels

$$y_{ic} = \mathbb{I}(c = y_i)$$

Multi-class logistic regression model

Loss function: Negative log likelihood (NLL)

$$\ell(\mathbf{W}) = -\log L(\mathbf{W}) = -\sum_{i=1}^n \sum_{c=1}^C y_{ic} \log \mu_{ic}$$

which is also called the cross-entropy loss (CE loss) function, it measures the Kullback-Leibler (KL) divergence between two distributions μ_{ic} and y_{ic}

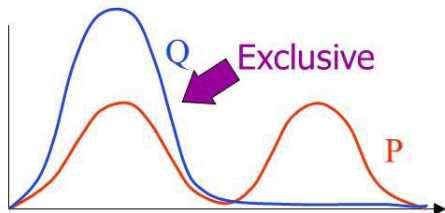
Recall: KL divergence between two distributions P and Q

$$D_{\text{KL}}(P\|Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$$

Exclusive versus Inclusive

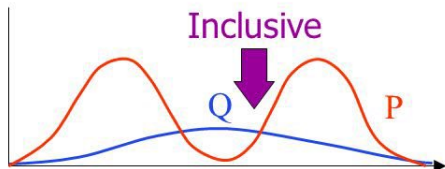
Minimising
 $KL(Q||P)$

$$= \sum_H Q(H) \ln \frac{Q(H)}{P(H|V)}$$



Minimising
 $KL(P||Q)$

$$= \sum_H P(H|V) \ln \frac{P(H|V)}{Q(H)}$$



Exclusive: Mode-seeking — Inclusive: Mean-seeking

Multi-class logistic regression model

Loss function: Negative log likelihood (NLL)

$$\ell(\mathbf{W}) = -\log L(\mathbf{W}) = -\sum_{i=1}^n \sum_{c=1}^C y_{ic} \log \mu_{ic}$$

which is also called the cross-entropy loss (CE loss) function, it measures the Kullback-Leibler (KL) divergence between two distributions μ_{ic} and y_{ic}

Question: What is the equivalent form of KL divergence of NLL, is it $\mathbb{E}[D_{\text{KL}}(y\|\mu)]$ or $\mathbb{E}[D_{\text{KL}}(\mu\|y)]$? Can we reverse the order? Why?

Multi-class logistic regression model

Gradient descend

The gradient¹ of the loss function w.r.t the loss function are

$$\nabla_{w_c} \ell(\mathbf{w}) = \sum_{i=1}^n (\mu_{ic} - y_{ic}) x_i$$

$$\nabla_{w_{c0}} \ell(\mathbf{w}) = \sum_{i=1}^n (\mu_{ic} - y_{ic})$$

¹<https://deepnotes.io/softmax-crossentropy>

Table of content

Multi-class logistics regression model

Multi-class classification

The optimal classifier

K nearest neighbor - KNN

Multi-class classification

Problem statement

Given a dataset $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y} = \{1, 2, \dots, C\}$, we want to find a classifier $h(x) : \mathcal{X} \rightarrow \mathcal{Y}$

Statistical probability perspective

The dataset D is sampled from an unknown distribution $\mathcal{P}(x, y)$.
A 'good' classifier h has a low probability of making error under \mathcal{P} .
Given a sample $(X, Y) \sim \mathcal{P}$, the event X is misclassified by h :

$$\{h(X) \neq Y\}$$

and the *error probability (error rate)* of h

$$\text{err}_{\mathcal{P}}(h) = \mathbb{P}_{X, Y \sim \mathcal{P}}\{h(X) \neq Y\} \quad (1)$$

Multi-class classification

Estimate the error rate

Since P is unknown, the true error (1) is not directly available. However, we have the dataset D as a realization of P . The empirical error rate on the training data D is then

$$\widehat{\text{err}}_D(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(h(x_i) \neq y_i) \quad (2)$$

Expectation of empirical error rate

The empirical error rate (2) is an unbiased estimator of the true error rate (1) under P since

$$\begin{aligned} \mathbb{E}_P[\widehat{\text{err}}_D(h)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{x_i, y_i \sim \mathcal{P}}[\mathbb{I}(h(x_i) \neq y_i)] \\ &= \mathbb{P}\{h(X) \neq Y\} = \text{err}_{\mathcal{P}}(h) \end{aligned}$$

Bounding Empirical Error Rate

Concentration bound²

Hoeffding inequality with empirical mean reminder: Let X_1, X_2, \dots, X_n be (random) samples from a random variable X , bounded a.s in $[a, b]$, then

$$\mathbb{P}\left\{\left|\frac{1}{n}\sum_i^n X_i - E[X]\right| \leq \epsilon\right\} \geq 1 - 2\exp\frac{-2n\epsilon^2}{(b-a)^2}$$

Connection to confidence intervals: Bound the range of error ϵ with the level of significance α (change to have a wrong estimation) with the needed number of examples n

²<https://www.stat.cmu.edu/larry/=stat700/Lecture6.pdf>

Bounding Empirical Error Rate

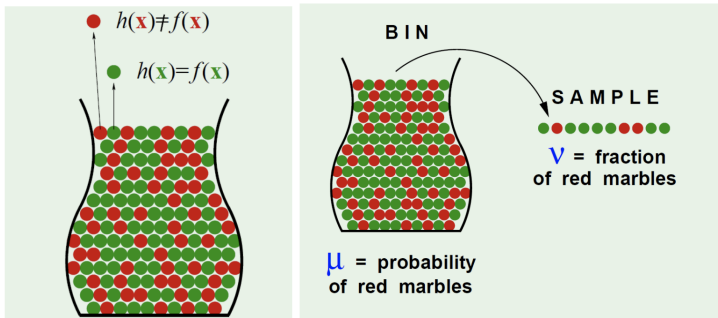
Concentration bound Apply the Hoeffding inequality to the empirical error rate, we have a concentration bound on the difference between empirical and true error rates,

$$\mathbb{P} \{ |\widehat{\text{err}}_D(h) - \text{err}_{\mathcal{P}}(h) | \leq \epsilon \} \geq 1 - 2e^{-2n\epsilon^2}$$

with $X_i = \mathbb{I}(h(x_i) \neq y_i)$ are indicator random variables taking values in $\{0, 1\}$.

- ▶ The inequality shows how $\widehat{\text{err}}_D(h)$ closed to $\text{err}_{\mathcal{P}}(h)$, given ϵ and the number of samples

Bounding Empirical Error Rate



- $X_n = \mathbb{I}[h(\mathbf{x}_n) \neq f(\mathbf{x}_n)] = \text{one sample training error} = \text{either 0 or 1}$

Table of content

Multi-class logistics regression model

Multi-class classification

The optimal classifier

K nearest neighbor - KNN

The optimal classifier

Fixed $X = x$, the error event

$$\mathbb{P}\{h(X) \neq Y | X = x\} = 1 - \mathbb{P}(Y = h(x) | X = x)$$

Question: If $X = x$ is fixed, what should be $h(x)$ from the numbers 1 to C to minimize the probability of error?

$$\begin{bmatrix} \mathbb{P}(Y = 1 | X = x) \\ \mathbb{P}(Y = 2 | X = x) \\ \vdots \\ \mathbb{P}(Y = C | X = x) \end{bmatrix}$$

Bayes optimal classifier: is a probabilistic model that makes the most probable classification

$$h^*(x) = \arg \max_c \mathbb{P}(Y = c | X = x)$$

The optimal classifier

Optimal error probability

$$E^* = \text{err}_{\mathcal{P}}(h^*) \leq \text{err}_{\mathcal{P}}(h), \forall h$$

For example: A logistic regression model uses the formula of the Bayes optimal classifier, with $\mathbb{P}(Y = c|X = x)$ is the probabilities of the categorical distribution $\text{Cat}(y|\mathcal{S}(f_1(x), f_2(x), \dots, f_C(x)))$

Question: Is it possible to find the Bayesian classification function?

- ▶ Can we find the distribution P ?
- ▶ Can we approximate $\mathbb{P}(Y = c|X = x)$ to approximate h^* ?

The optimal classifier

Two approaches in machine learning

- ▶ **Discriminative models:** learn a predictor given the observations
 - ▶ Approximate $\mathbb{P}(Y = c|X = x)$
 - ▶ Approximate the classifier h^*
- ▶ **Generative models:** learn a joint distribution over all the variables.
 - ▶ Approximate $\mathbb{P}(Y = c, X = x)$

$$\begin{aligned}h^*(x) &= \arg \max_c \mathbb{P}(Y = c, X = x) \\&= \arg \max_c \mathbb{P}(X = x|Y = c)\mathbb{P}(Y = c)\end{aligned}$$

- ▶ Can solve a variety of problems, not only classification

Discriminative vs generative models

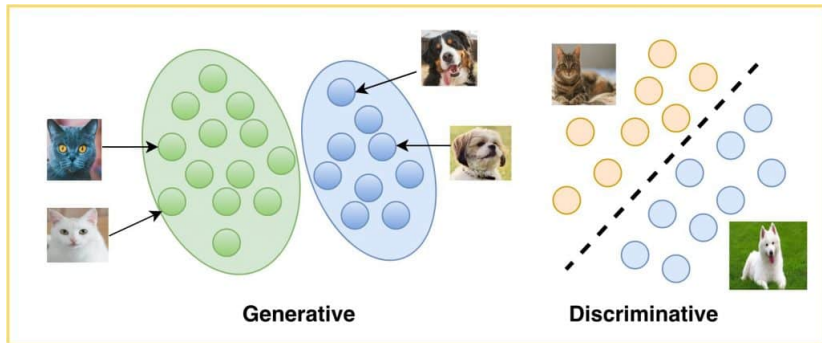


Table of content

Multi-class logistics regression model

Multi-class classification

The optimal classifier

K nearest neighbor - KNN

K nearest neighbor - KNN

Estimate $\mathbb{P}(Y = c|X = x)$ in the neighborhood of x

- ▶ Find k samples from D that are closest to x
- ▶ Approximate $\mathbb{P}(Y = c|X = x)$ by the proportion of label c in the k samples.

$$\hat{\mathbb{P}}(Y = c|X = x) = \frac{\sum_{i=1}^k \mathbb{I}(x_i \in V_k(x)) \mathbb{I}(y_i = c)}{k}$$

where $V_k(x)$ is a neighborhood of x containing k of data samples in D . The numerator is the number of data samples in $V_k(x)$ that are labeled c .

$h^{KNN}(x)$: select the label c that occurs most in k of the data sample closest to x in D .

K nearest neighbor - KNN

Pseudocode of KNN

1. Calculate $d(x, x_i)$, $i = 1, 2, \dots, n$, where d is a distance metric.
2. Take the first k data points whose distances are smallest among the calculated distance list, along with their labels, denoted $(x_j, y_j)_{j=1}^k$.
3. Return the label which has the most votes.

K nearest neighbor - KNN

Theorem of the upper bound of error rate of knn with $k = 1$

When the number of training samples $n \rightarrow \infty$, we have

$$R^* \leq \text{err}_{\mathcal{P}}(h^{\text{KNN}}) \leq R^* \left(2 - \frac{C}{C-1} R^* \right).$$

where R^* is the Bayesian optimal error probability.

K nearest neighbor - KNN

Proof: Let $x_{(1)}$ as the nearest neighbor of x in D and $y_{(1)}$ as the label of this data sample and y^{true} as the label of x .

Suppose that when number of samples $n \rightarrow \infty$,

$$\begin{aligned}x_{(1)} &\rightarrow x \\ P(y|x_{(1)}) &\rightarrow P(y|x), \forall y = 1, 2, \dots, C\end{aligned}$$

The error probability of h^{KNN} on x occurs when $y_{(1)} \neq y^{true}$ is:

$$\begin{aligned}\text{err}(h^{KNN}, x) &= P(y^{true} \neq y_{(1)} | x, x_{(1)}) \\ &= 1 - \sum_{y=1}^C P(y^{true} = y | x) P(y_{(1)} = y | x_{(1)}) \\ &\rightarrow 1 - \sum_{y=1}^C P^2(y | x) \quad (\text{as } n \rightarrow \infty)\end{aligned}$$

K nearest neighbor - KNN

If $y^* = h^*(x)$ is the output of the Bayesian optimal classifier function, then

$$\begin{aligned}P(y^*|x) &= \max_y P(y|x) = 1 - \text{err}(h^*, x) \\ \sum_{y=1}^C P^2(y|x) &= P^2(y^*|x) + \sum_{y \neq y^*} P^2(y|x) \\ &\geq P^2(y^*|x) + \frac{(\sum_{y \neq y^*} P(y|x))^2}{C-1} \\ &\quad \text{(Bunyakovsky inequality)} \\ &= (1 - \text{err}(h^*, x))^2 + \frac{\text{err}(h^*, x)^2}{C-1}\end{aligned}$$

So

$$1 - \sum_{y=1}^C P^2(y|x) \leq 2\text{err}(h^*, x) - \frac{C}{C-1}\text{err}(h^*, x)^2 \quad (3)$$

K nearest neighbor - KNN

Taking the expectation both sides of (3), with $R^* = \mathbb{E}_x \text{err}(h^*, x)$, we have

$$\mathbb{E}_x \text{err}(h^{\text{KNN}}, x) = \text{err}(h^{\text{KNN}}) \leq 2R^* - \frac{C}{C-1} \int_x \text{err}(h^*, x)^2 p(x) dx \quad (4)$$

Also,

$$\begin{aligned} 0 &\leq \int_x (\text{err}(h^*, x) - R^*)^2 p(x) dx \\ &= \int_x (\text{err}(h^*, x)^2 - 2R^* \text{err}(h^*, x) + (R^*)^2) p(x) dx \\ &= \int_x \text{err}(h^*, x)^2 p(x) dx - (R^*)^2 \end{aligned}$$

Substitute to (4)

$$\text{err}(h^{\text{KNN}}) \leq 2R^* - \frac{C}{C-1} (R^*)^2 = R^* \left[2 - \frac{C}{C-1} R^* \right].$$