Machine learning - HT 2016 6. Classification: Logistic Regression

Varun Kanade

University of Oxford February 10, 2016

Outline

Today we'll discuss classification using logistic regression.

- Discriminative vs Generative Models
- Likelihood of Logistic Regression
- Using convex optimization to the obtain MLE
- Logistic Regression in torch

Classification : Generative Models

How are the inputs, tail length and height, distributed given the class?

Model $Pr(\mathbf{x} \mid y = \mathsf{zebra})$

Model $Pr(\mathbf{x} \mid y = \mathsf{donkey})$

Example: Model both distributions are multivariate normal with same covariance matrix but different mean



Classification : Generative Models



Don't try to model the inputs ${\bf x}$ at all

Model the output y given the input \mathbf{x} and the parameters for the model \mathbf{w}

 $y \sim p(\mathbf{x}, \mathbf{w})$

Don't try to model the inputs **x** at all

Model the output y given the input \mathbf{x} and the parameters for the model \mathbf{w}

 $y \sim p(\mathbf{x}, \mathbf{w})$

Pros and cons for both approaches (see Murphy Chapter 8.6)

Focus on discriminative classification

Logistic Regression: Sigmoid Function

The sigmoid function, or σ , (a.k.a. logistic or logit) is defined as

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



Binary Classification: Logistic Regression

As in the case of linear regression, we model y given $\mathbf{x}\in\mathbb{R}^n$ and parameters $\mathbf{w}\in\mathbb{R}^n$

Linear model parametrized by $\mathbf{w} \in \mathbb{R}^n$ composed with sigmoid filter

We have,

$$\Pr(y = 1 \mid \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$$

Binary Classification: Logistic Regression

As in the case of linear regression, we model y given $\mathbf{x}\in\mathbb{R}^n$ and parameters $\mathbf{w}\in\mathbb{R}^n$

Linear model parametrized by $\mathbf{w} \in \mathbb{R}^n$ composed with sigmoid filter

We have,

$$\Pr(y = 1 \mid \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{x}^T \mathbf{w})$$

For prediction:

$$\hat{y} = \mathbb{I}(\sigma(\mathbf{x}^T \mathbf{w}) \ge \frac{1}{2})$$

Binary Classification : Logistic Regression





Bernoulli Random Variables

Bernoulli random variable X takes value in $\{0,1\}.$ We parametrize using $\theta \in [0,1].$

$$p(1 \mid \theta) = \theta$$
$$p(0 \mid \theta) = 1 - \theta$$

Bernoulli Random Variables

Bernoulli random variable X takes value in $\{0,1\}.$ We parametrize using $\theta \in [0,1].$

$$p(1 \mid \theta) = \theta$$
$$p(0 \mid \theta) = 1 - \theta$$

More succinctly, we can write

$$p(x \mid \theta) = \theta^{x} (1 - \theta)^{1 - x}$$

Bernoulli Random Variables

Bernoulli random variable X takes value in $\{0, 1\}$. We parametrize using $\theta \in [0, 1]$.

$$p(1 \mid \theta) = \theta$$
$$p(0 \mid \theta) = 1 - \theta$$

More succinctly, we can write

$$p(x \mid \theta) = \theta^{x} (1 - \theta)^{1 - x}$$

Logistic Regression

y given x and parameter w is modelled as Bernoulli variable

 $y \sim \text{Bernoulli}(\sigma(\mathbf{x}^T \mathbf{w}))$

Likelihood of Logistic Regression

Given data $D = \langle (\mathbf{x}_i, y_i) \rangle_{i=1}^m$ we can compute the likelihood of observing \mathbf{y} under the logistic regression model

Likelihood of Logistic Regression

Given data $D = \langle (\mathbf{x}_i, y_i) \rangle_{i=1}^m$ we can compute the likelihood of observing \mathbf{y} under the logistic regression model

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \prod_{i=1}^{m} \text{Bernoulli}(y_i \mid \sigma(\mathbf{x}_i^T \mathbf{w}))$$
$$= \prod_{i=1}^{m} \left(\frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}\right)^{y_i} \left(1 - \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}\right)^{1 - y_i}$$

Likelihood of Logistic Regression

Given data $\mathcal{D} = \langle (\mathbf{x}_i, y_i) \rangle_{i=1}^m$ we can compute the likelihood of observing \mathbf{y} under the logistic regression model

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \prod_{i=1}^{m} \operatorname{Bernoulli}(y_i \mid \sigma(\mathbf{x}_i^T \mathbf{w}))$$
$$= \prod_{i=1}^{m} \left(\frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}\right)^{y_i} \left(1 - \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}}\right)^{1 - y_i}$$

Let's look at the negative log likelihood for a single data point (\mathbf{x}_i, y_i)

$$L(\mathbf{w}; \mathbf{x}_i, y_i) = -\log(p(y_i \mid \sigma(\mathbf{x}_i^T \mathbf{w})))$$

= - (y_i log(\pi_i) + (1 - y_i) log(1 - \pi_i))

Gradient and Hessian of NLL

The negative log likelihood is given by

$$L(\mathbf{w}) = \text{NLL}(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = -\sum_{i=1}^{m} (y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i))$$

Gradient and Hessian of NLL

The negative log likelihood is given by

$$L(\mathbf{w}) = \text{NLL}(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = -\sum_{i=1}^{m} (y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i))$$

The gradient and the Hessian (with respect to w) can be computed as:

$$\mathbf{g} = \nabla_{\mathbf{w}} \mathbf{L} = \sum_{i=1}^{m} \mathbf{x}_i (\pi_i - y_i) = \mathbf{X}^T (\boldsymbol{\pi} - \mathbf{y})$$
$$\mathbf{H} = \nabla_{\mathbf{w}}^2 \mathbf{L} = \sum_{i=1}^{m} \pi_i (1 - \pi_i) \mathbf{x}^i \mathbf{x}_i^T = \mathbf{X}^T \operatorname{diag}(\pi_i (1 - \pi_i)) \mathbf{X}$$

Gradient and Hessian of NLL

The negative log likelihood is given by

$$L(\mathbf{w}) = \text{NLL}(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = -\sum_{i=1}^{m} (y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i))$$

The gradient and the Hessian (with respect to \mathbf{w}) can be computed as:

$$\mathbf{g} = \nabla_{\mathbf{w}} \mathbf{L} = \sum_{i=1}^{m} \mathbf{x}_{i}(\pi_{i} - y_{i}) = \mathbf{X}^{T}(\boldsymbol{\pi} - \mathbf{y})$$
$$\mathbf{H} = \nabla_{\mathbf{w}}^{2} \mathbf{L} = \sum_{i=1}^{m} \pi_{i}(1 - \pi_{i})\mathbf{x}^{i}\mathbf{x}_{i}^{T} = \mathbf{X}^{T} \operatorname{diag}(\pi_{i}(1 - \pi_{i}))\mathbf{X}$$

Homework: Show that H is positive definite.

NLL is convex and has a global minimum

Iteratively Reweighted Least Squares (IRLS)

Apply Newton's method

$$\begin{aligned} \mathbf{g}_t &= \mathbf{X}^T (\boldsymbol{\pi}_t - \mathbf{y}) = -\mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}_t) \\ \mathbf{H}_t &= \mathbf{X}^T \mathbf{S}_t \mathbf{X} \end{aligned}$$

Iteratively Reweighted Least Squares (IRLS)

Apply Newton's method

$$\mathbf{g}_t = \mathbf{X}^T (\boldsymbol{\pi}_t - \mathbf{y}) = -\mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}_t)$$

 $\mathbf{H}_t = \mathbf{X}^T \mathbf{S}_t \mathbf{X}$

Newton's update says:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \mathbf{H}_t^{-1} \mathbf{g}_t \\ &= \mathbf{w}_t + (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}_t) \\ &= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{S}_t \mathbf{X} \mathbf{w}_t + \mathbf{y} - \boldsymbol{\pi}_t) = (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} (\mathbf{X}^t \mathbf{S}_t \mathbf{z}_t) \end{aligned}$$

Iteratively Reweighted Least Squares (IRLS)

Apply Newton's method

$$\mathbf{g}_t = \mathbf{X}^T (\boldsymbol{\pi}_t - \mathbf{y}) = -\mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}_t)$$

 $\mathbf{H}_t = \mathbf{X}^T \mathbf{S}_t \mathbf{X}$

Newton's update says:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \mathbf{H}_t^{-1} \mathbf{g}_t \\ &= \mathbf{w}_t + (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}_t) \\ &= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{S}_t \mathbf{X} \mathbf{w}_t + \mathbf{y} - \boldsymbol{\pi}_t) = (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} (\mathbf{X}^t \mathbf{S}_t \mathbf{z}_t) \end{aligned}$$

This is a least square solution for the system

$$\sum_{i=1}^{m} S_{t,i} (\mathbf{x}_i^T \mathbf{w} - z_{t,i})^2$$

Multi-class, Softmax Formulation, Multinoulli¹

Logistic Regression as a Neural Network

¹Kevin Murphy's usage

Softmax in Torch

nn.SoftMax()



Likelihood for multi-class

Classes: $\{1, ..., C\}$

Indicator function:

$$\mathbb{I}_c(y) = egin{cases} 1 & \text{if } y = c \ 0 & \text{otherwise} \end{cases}$$

The parameters ${f W}$ is now a n imes C matrix

For a single data point (x, y) the likelihood is:

$$p(y \mid \mathbf{x}, \mathbf{W}) = \prod_{c=1}^{C} \pi_{c}^{\mathbb{I}_{c}(y)}$$

And the negative log likelihood is

$$L(\mathbf{W}; \mathbf{x}, \mathbf{y}) = -\sum_{c=1}^{C} \mathbb{I}_{c}(y) \log(\pi_{c})$$

Multiclass Logistic Regression in Torch

```
example-logistic-regression.lua
require 'nn'; require 'optim';
model = nn.Sequential()
ninputs = 10; noutputs = 3
model:add(nn.Linear(ninputs, noutputs))
model:add(nn.LogSoftMax())
criterion = nn.ClassNLLCriterion()
-- define some input and target
-- to evaluate model
model:forward(input)
-- to evaluate loss
criterion:forward(model:forward(input), target)
-- to compute gradients
model:backward(input, criterion:backward(model.output, target))
```