Logistic regression

Daniel Hsu (COMS 4771)

The logistic regression model

Logistic regression is a model for binary classification data with feature vectors in \mathbb{R}^d and labels in $\{-1, +1\}$. Data $(\boldsymbol{X}_1, Y_1), \dots, (\boldsymbol{X}_n, Y_n)$ are treated as iid random variables taking values in $\mathbb{R}^d \times \{-1, +1\}$, and for each $\boldsymbol{x} \in \mathbb{R}^d$.

$$Y_i \mid \boldsymbol{X}_i = \boldsymbol{x} \sim \operatorname{Bern}(\sigma(\boldsymbol{x}^{\mathsf{T}} \boldsymbol{w}))$$

where $\sigma(t) = 1/(1 + \exp(-t))$ is the sigmoid function. Here, $\mathbf{w} \in \mathbb{R}^d$ is the parameter of the model, and it is not involved in the marginal distribution of \mathbf{X}_i (which we leave unspecified).

Maximum likelihood

The log-likelia Assignment, Project, Exam Help

w)

There is no closed-form that ps://eduassistpro.github.io/mately minimize the negative log-li

Empirical risk Andriz WorChat edu_assist_pro

Maximum likelihood is very different from finding the linear classifier of smallest empirical zero-one loss risk. Finding the empirical zero-one loss risk minimizer is computationally intractable in general.

Finding a linear separator

There are special cases when finding the empirical zero-one loss risk minimizer is computationally tractable. One is when the training data is *linearly separable*: i.e., when there exists $\boldsymbol{w}^{\star} \in \mathbb{R}^{d}$ such that

$$y_i \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{w}^{\star} > 0$$
, for all $i = 1, \dots, n$.

Claim. Define $L(\boldsymbol{w}) := \sum_{i=1}^n \ln(1 + \exp(-y_i \boldsymbol{x}_i^{\mathsf{T}} \boldsymbol{w}))$. Suppose $(\boldsymbol{x}_1, y_1), \dots, (\boldsymbol{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ is linearly separable. Then any $\hat{\boldsymbol{w}} \in \mathbb{R}^d$ with

$$L(\hat{\boldsymbol{w}}) < \inf_{\boldsymbol{w} \in \mathbb{R}^d} L(\boldsymbol{w}) + \ln(2)$$

is a linear separator.

Proof. We first observe that the infimum¹ (i.e., greatest lower bound) of L is zero. Let $\mathbf{w}^* \in \mathbb{R}^d$ be a linear separator, so $s_i := y_i \mathbf{x}_i^{\mathsf{T}} \mathbf{w}^* > 0$ for all $i = 1, \ldots, n$. For any r > 0,

$$L(r\boldsymbol{w}^{\star}) = \sum_{i=1}^{n} \ln(1 + \exp(-rs_i)),$$

¹https://en.wikipedia.org/wiki/Infimum_and_supremum

and therefore

$$\lim_{r \to \infty} \sum_{i=1}^{n} \ln(1 + \exp(-rs_i)) = 0.$$

Every term $\ln(1 + \exp(-y_i \boldsymbol{x}_i^{\dagger} \boldsymbol{w}))$ in $L(\boldsymbol{w})$ is positive, so $L(\boldsymbol{w}) > 0$. Therefore, we conclude that

$$\inf_{\boldsymbol{w}\in\mathbb{R}^d}L(\boldsymbol{w})=0.$$

So now we just have to show that any $\hat{\boldsymbol{w}} \in \mathbb{R}^d$ with

$$L(\hat{\boldsymbol{w}}) < \ln(2)$$

is a linear separator. So let $\hat{\boldsymbol{w}}$ satisfy $L(\hat{\boldsymbol{w}}) < \ln(2)$, which implies

$$\ln(1 + \exp(-y_i \boldsymbol{x}_i^{\mathsf{T}} \hat{\boldsymbol{w}})) < \ln(2)$$

for every i = 1, ..., n. Exponentiating both sides gives

$$1 + \exp(-y_i \boldsymbol{x}_i^{\mathsf{T}} \hat{\boldsymbol{w}}) < 2.$$

Now subtracting 1 from both sides and taking logarithms gives

Assignment Project Exam Help This means that $\hat{\boldsymbol{w}}$ correctly classifies (\boldsymbol{x}_i, y_i) . Since this holds for all i = 1, ..., n, it follows that $\hat{\boldsymbol{w}}$ is a

Surrogate loss https://eduassistpro.github.io/

Even if $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ is not linearles imizing the log-likelihood can yield a problem of the log-likelihood can yield a p

$$\widehat{\mathcal{R}}(\boldsymbol{w}) := \frac{1}{n} \sum_{i=1}^{n} \ell_{\log}(y_i \boldsymbol{x}_i^{\scriptscriptstyle\mathsf{T}} \boldsymbol{w})$$

where

$$\ell_{\log}(z) := -\ln \sigma(z)$$

is the logistic loss. The logistic loss (up to scaling) turns out to be an upper-bound on the zero-one loss:

$$\ell_{\mathrm{zo}}(z) \le \frac{1}{\ln 2} \, \ell_{\mathrm{log}}(z),$$

where $\ell_{zo}(z) = \mathbb{1}_{\{z \le 0\}}$. If the empirical logistic loss risk is small, then the empirical zero-one loss is also small.

Gradient descent for logistic regression

The derivative of ℓ_{\log} is given by

$$\frac{\mathrm{d}\ell_{\log}(z)}{\mathrm{d}z} = -\frac{1}{\sigma(z)} \cdot \frac{\mathrm{d}\sigma(z)}{\mathrm{d}z}$$
$$= -\frac{1}{\sigma(z)} \cdot \sigma(z) \cdot \sigma(-z)$$
$$= -\sigma(-z).$$

Therefore, by linearity and the chain rule, the negative gradient of $\widehat{\mathcal{R}}$ with respect to \boldsymbol{w} is

$$\begin{split} -\nabla \widehat{\mathcal{R}}(\boldsymbol{w}) &= -\frac{1}{n} \sum_{i=1}^{n} \nabla \ell_{\log}(y_{i} \boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{w}) \\ &= -\frac{1}{n} \sum_{i=1}^{n} \frac{\mathrm{d}\ell_{\log}(z)}{\mathrm{d}z} \bigg|_{z=y_{i} \boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{w}} \cdot \nabla \left(y_{i} \boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{w}\right) \\ &= \frac{1}{n} \sum_{i=1}^{n} \sigma(-y_{i} \boldsymbol{x}_{i}^{\mathsf{T}} \boldsymbol{w}) \cdot y_{i} \boldsymbol{x}_{i}. \end{split}$$

Now suppose $\mathbf{A} = [\mathbf{x}_1|\cdots|\mathbf{x}_n]^{\mathsf{T}} \in \mathbb{R}^{n \times d}$ and $\mathbf{b} = [y_1|\cdots|y_n]^{\mathsf{T}} \in \mathbb{R}^n$. (Notice that we have omitted the $1/\sqrt{n}$ scaling that we had for least squares linear regression.) Then the negative gradient of $\widehat{\mathcal{R}}$ can be written as

$$-\nabla \widehat{\mathcal{R}}(\boldsymbol{w}) = \frac{1}{n} \boldsymbol{A}^{\mathsf{T}} (\boldsymbol{b} \odot \sigma(-\boldsymbol{b} \odot (\boldsymbol{A}\boldsymbol{w}))),$$

where $\boldsymbol{u} \odot \boldsymbol{v} \in \mathbb{R}^n$ is the coordinate-wise product of vectors $\boldsymbol{u} \in \mathbb{R}^n$ and $\boldsymbol{v} \in \mathbb{R}^n$, and $\sigma(\boldsymbol{v}) \in \mathbb{R}^n$ is the coordinate-wise application of the sigmoid function to $\boldsymbol{v} \in \mathbb{R}^n$.

Gradient descent for logistic regression begins with an initial weight vector $\boldsymbol{w}^{(0)} \in \mathbb{R}^d$, and then iteratively updates it by subtracting a positive multiple $\eta > 0$ of the gradient at the current iterate:

Assignment-Project Exam Help

https://eduassistpro.github.io/ Add WeChat edu_assist_pro