

Introduction to Machine Learning

Logistic Regression

Varun Chandola

Computer Science & Engineering
State University of New York at Buffalo
Buffalo, NY, USA
chandola@buffalo.edu



University at Buffalo
Department of Computer Science
and Engineering
School of Engineering and Applied Sciences

Generative vs. Discriminative Classifiers

Logistic Regression

Logistic Regression - Training

- Using Gradient Descent for Learning Weights

- Using Newton's Method

- Regularization with Logistic Regression

- Handling Multiple Classes

- Bayesian Logistic Regression

- Laplace Approximation

- Posterior of \mathbf{w} for Logistic Regression

- Approximating the Posterior

- Getting Prediction on Unseen Examples

Generative vs. Discriminative Classifiers

- ▶ Probabilistic classification task:

$$p(Y = \textit{benign} | \mathbf{X} = \mathbf{x}), p(Y = \textit{malicious} | \mathbf{X} = \mathbf{x})$$

- ▶ How do you estimate $p(y|\mathbf{x})$?

$$p(y|\mathbf{x}) = \frac{p(y, \mathbf{x})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})}$$

- ▶ Two step approach - Estimate generative model and then posterior for y (Naïve Bayes)
- ▶ Solving a more general problem [2, 1]
- ▶ Why not directly model $p(y|\mathbf{x})$? - **Discriminative approach**

Which is Better?

- ▶ Number of training examples needed to learn a PAC-learnable classifier \propto *VC-dimension of the hypothesis space*
- ▶ VC-dimension of a probabilistic classifier \propto Number of parameters [2] (or a small polynomial in the number of parameters)
- ▶ Number of parameters for $p(y, \mathbf{x}) >$ Number of parameters for $p(y|\mathbf{x})$

Discriminative classifiers need lesser training examples to for PAC learning than generative classifiers

Logistic Regression

- ▶ $y|\mathbf{x}$ is a *Bernoulli* distribution with parameter $\theta = \text{sigmoid}(\mathbf{w}^\top \mathbf{x})$
- ▶ When a new input \mathbf{x}^* arrives, we toss a coin which has $\text{sigmoid}(\mathbf{w}^\top \mathbf{x}^*)$ as the probability of heads
- ▶ If outcome is heads, the predicted class is 1 else 0
- ▶ Learns a linear boundary

Learning Task for Logistic Regression

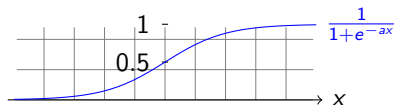
Given training examples $\langle \mathbf{x}_i, y_i \rangle_{i=1}^D$, learn \mathbf{w}

Bayesian Interpretation

- ▶ Directly model $p(y|\mathbf{x})$ ($y \in \{0, 1\}$)
- ▶ $p(y|\mathbf{x}) \sim \text{Bernoulli}(\theta = \text{sigmoid}(\mathbf{w}^\top \mathbf{x}))$

Geometric Interpretation

- ▶ Use regression to predict discrete values
- ▶ *Squash* output to $[0, 1]$ using sigmoid function
- ▶ Output less than 0.5 is one class and greater than 0.5 is the other



- ▶ MLE Approach
- ▶ Assume that $y \in \{0, 1\}$
- ▶ What is the likelihood for a bernoulli sample?
 - ▶ If $y_i = 1$, $p(y_i) = \theta_i = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_i)}$
 - ▶ If $y_i = 0$, $p(y_i) = 1 - \theta_i = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x}_i)}$
 - ▶ In general, $p(y_i) = \theta_i^{y_i} (1 - \theta_i)^{1 - y_i}$

Log-likelihood

$$LL(\mathbf{w}) = \sum_{i=1}^N y_i \log \theta_i + (1 - y_i) \log (1 - \theta_i)$$

- ▶ No closed form solution for maximizing log-likelihood

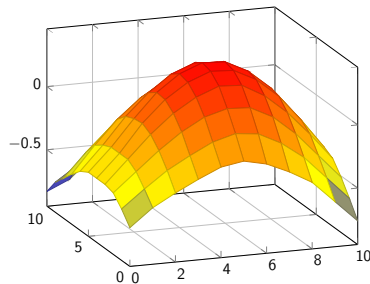
Using Gradient Descent for Learning Weights

- ▶ Compute gradient of LL with respect to \mathbf{w}
- ▶ A convex function of \mathbf{w} with a unique global maximum

$$\frac{d}{d\mathbf{w}} LL(\mathbf{w}) = \sum_{i=1}^N (y_i - \theta_i) \mathbf{x}_i$$

- ▶ Update rule:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta \frac{d}{d\mathbf{w}_k} LL(\mathbf{w}_k)$$



Using Newton's Method

- ▶ Setting η is sometimes *tricky*
- ▶ Too large – incorrect results
- ▶ Too small – slow convergence
- ▶ Another way to speed up convergence:

Newton's Method

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta \mathbf{H}_k^{-1} \frac{d}{d\mathbf{w}_k} LL(\mathbf{w}_k)$$

What is the Hessian?

- ▶ Hessian or **H** is the second order derivative of the objective function
- ▶ Newton's method belong to the family of **second order optimization algorithms**
- ▶ For logistic regression, the Hessian is:

$$H = - \sum_i \theta_i(1 - \theta_i) \mathbf{x}_i \mathbf{x}_i^\top$$

Regularization with Logistic Regression

- ▶ **Overfitting** is an issue, especially with large number of features
- ▶ Add a *Gaussian prior* $\sim \mathcal{N}(\mathbf{0}, \tau^2)$ (Or a regularization penalty)
- ▶ Easy to incorporate in the gradient descent based approach

$$LL'(\mathbf{w}) = LL(\mathbf{w}) - \frac{1}{2}\lambda\mathbf{w}^\top\mathbf{w}$$

$$\frac{d}{d\mathbf{w}}LL'(\mathbf{w}) = \frac{d}{d\mathbf{w}}LL(\mathbf{w}) - \lambda\mathbf{w}$$

$$H' = H - \lambda I$$

where I is the identity matrix.

Handling Multiple Classes

- ▶ One vs. Rest and One vs. Other
- ▶ $p(y|\mathbf{x}) \sim \text{Multinoulli}(\boldsymbol{\theta})$
- ▶ Multinoulli parameter vector $\boldsymbol{\theta}$ is defined as:

$$\theta_j = \frac{\exp(\mathbf{w}_j^\top \mathbf{x})}{\sum_{k=1}^C \exp(\mathbf{w}_k^\top \mathbf{x})}$$

- ▶ Multiclass logistic regression has C weight vectors to learn

- ▶ How to get the posterior for \mathbf{w} ?
- ▶ Not easy - *Why?*

Laplace Approximation

- ▶ We do not know what the true posterior distribution for \mathbf{w} is.
- ▶ Is there a close-enough (approximate) Gaussian distribution?

Problem Statement

How to approximate a posterior with a Gaussian distribution?

- ▶ When is this needed?
 - ▶ When direct computation of posterior is not possible.
 - ▶ No conjugate prior ☹

Laplace Approximation using Taylor Series Expansion

Laplace Approximation using Taylor Series Expansion

- ▶ Assume that posterior is:

$$p(\mathbf{w}|D) = \frac{1}{Z} e^{-E(\mathbf{w})}$$

$E(\mathbf{w})$ is **energy function** \equiv negative log of unnormalized log posterior

Laplace Approximation using Taylor Series Expansion

- ▶ Assume that posterior is:

$$p(\mathbf{w}|D) = \frac{1}{Z} e^{-E(\mathbf{w})}$$

$E(\mathbf{w})$ is **energy function** \equiv negative log of unnormalized log posterior

- ▶ Let \mathbf{w}_{MAP} be the *mode* or *expected value* of the posterior distribution of \mathbf{w}

Laplace Approximation using Taylor Series Expansion

- ▶ Assume that posterior is:

$$p(\mathbf{w}|D) = \frac{1}{Z} e^{-E(\mathbf{w})}$$

$E(\mathbf{w})$ is **energy function** \equiv negative log of unnormalized log posterior

- ▶ Let \mathbf{w}_{MAP} be the *mode* or *expected value* of the posterior distribution of \mathbf{w}
- ▶ Taylor series expansion of $E(\mathbf{w})$ around the mode

$$\begin{aligned} E(\mathbf{w}) &= E(\mathbf{w}_{MAP}) + (\mathbf{w} - \mathbf{w}^*)^\top E'(\mathbf{w}) + (\mathbf{w} - \mathbf{w}^*)^\top E''(\mathbf{w})(\mathbf{w} - \mathbf{w}^*) + \dots \\ &\approx E(\mathbf{w}_{MAP}) + (\mathbf{w} - \mathbf{w}^*)^\top E'(\mathbf{w}) + (\mathbf{w} - \mathbf{w}_{MAP})^\top E''(\mathbf{w})(\mathbf{w} - \mathbf{w}^*) \end{aligned}$$

$E'(\mathbf{w}) = \nabla$ - first derivative of $E(\mathbf{w})$ (**gradient**) and $E''(\mathbf{w}) = \mathbf{H}$ is the second derivative (**Hessian**)

Taylor Series Expansion Continued

- ▶ Since \mathbf{w}_{MAP} is the mode, the first derivative or gradient is zero

$$E(\mathbf{w}) \approx E(\mathbf{w}_{MAP}) + (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

Taylor Series Expansion Continued

- ▶ Since \mathbf{w}_{MAP} is the mode, the first derivative or gradient is zero

$$E(\mathbf{w}) \approx E(\mathbf{w}_{MAP}) + (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

- ▶ Posterior $p(\mathbf{w}|D)$ may be written as:

$$\begin{aligned} p(\mathbf{w}|D) &\approx \frac{1}{Z} e^{-E(\mathbf{w}_{MAP})} \exp \left[-\frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*) \right] \\ &= \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{H}^{-1}) \end{aligned}$$

\mathbf{w}_{MAP} is the mode obtained by maximizing the posterior using gradient ascent

Posterior of \mathbf{w} for Logistic Regression

- ▶ Prior:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I})$$

- ▶ Likelihood of data

$$p(D|\mathbf{w}) = \prod_{i=1}^N \theta_i^{y_i} (1 - \theta_i)^{1-y_i}$$

where $\theta_i = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}_i}}$

- ▶ Posterior:

$$p(\mathbf{w}|D) = \frac{\mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}) \prod_{i=1}^N \theta_i^{y_i} (1 - \theta_i)^{1-y_i}}{\int p(D|\mathbf{w}) d\mathbf{w}}$$

Approximating the Posterior - Laplace Approximation

- ▶ Approximate posterior distribution

$$p(\mathbf{w}|D) = \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{H}^{-1})$$

- ▶ \mathbf{H} is the *Hessian* of the negative log-posterior w.r.t. \mathbf{w}

$$p(y|\mathbf{x}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|D)d\mathbf{w}$$

1. Use a point estimate of \mathbf{w} (MLE or MAP)
2. Analytical Result
3. **Monte Carlo Approximation**
 - ▶ Numerical integration
 - ▶ Sample finite "versions" of \mathbf{w} using $p(\mathbf{w}|D)$

$$p(\mathbf{w}|D) = \mathcal{N}(\mathbf{w}_{MAP}, \mathbf{H}^{-1})$$

- ▶ Compute $p(y|\mathbf{x})$ using the samples and add



A. Y. Ng and M. I. Jordan.

On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes.

In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *NIPS*, pages 841–848. MIT Press, 2001.



V. Vapnik.

Statistical learning theory.

Wiley, 1998.