

Probabilistic Linear Classification: Logistic Regression

Piyush Rai
IIT Kanpur

Probabilistic Machine Learning (CS772A)

Jan 18, 2016

Probabilistic Classification

- Given: N labeled training examples $\{\mathbf{x}_n, y_n\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \{0, 1\}$
- \mathbf{X} : $N \times D$ feature matrix, \mathbf{y} : $N \times 1$ label vector
- $y_n = 1$: positive example, $y_n = 0$: negative example
- Goal: Learn a classifier that predicts the binary label y_* for a new input \mathbf{x}_*
- Want a **probabilistic model** to be able to also predict the *label probabilities*

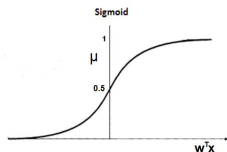
$$\begin{aligned}p(y_n = 1 | \mathbf{x}_n, \mathbf{w}) &= \mu_n \\p(y_n = 0 | \mathbf{x}_n, \mathbf{w}) &= 1 - \mu_n\end{aligned}$$

- $\mu_n \in (0, 1)$ is the probability of y_n being 1
- Note: Features \mathbf{x}_n assumed fixed (given). Only labels y_n being modeled
- \mathbf{w} is the model parameter (to be learned)
- How do we define μ_n (want it to be a function of \mathbf{w} and input \mathbf{x}_n)?

Logistic Regression

- Logistic regression defines μ using the **sigmoid function**

$$\mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



- Sigmoid computes a real-valued “score” ($\mathbf{w}^\top \mathbf{x}$) for input \mathbf{x} and “squashes” it between (0,1) to turn this score into a **probability** (of \mathbf{x} 's label being 1)
- Thus we have

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

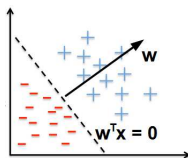
- Note:** If we assume $y \in \{-1, +1\}$ instead of $y \in \{0, 1\}$ then $p(y | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-y \mathbf{w}^\top \mathbf{x})}$

Logistic Regression: A Closer Look..

- What's the underlying decision rule in Logistic Regression?
- At the decision boundary, both classes are equiprobable. Thus:

$$\begin{aligned}p(y = 1|x, \mathbf{w}) &= p(y = 0|x, \mathbf{w}) \\ \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} &= \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \\ \exp(\mathbf{w}^\top \mathbf{x}) &= 1 \\ \mathbf{w}^\top \mathbf{x} &= 0\end{aligned}$$

- Thus the decision boundary of LR is nothing but a **linear hyperplane**, just like Perceptron, Support Vector Machine (SVM), etc.
- Therefore $y = 1$ if $\mathbf{w}^\top \mathbf{x} \geq 0$, otherwise $y = 0$

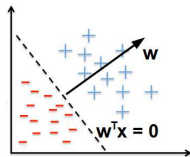
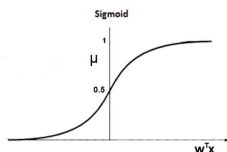


Interpreting the probabilities..

- Recall that

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = \mu = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

- Note that the “score” $\mathbf{w}^\top \mathbf{x}$ is also a measure of distance of \mathbf{x} from the hyperplane (score is positive for pos. examples, negative for neg. examples)



- High positive score $\mathbf{w}^\top \mathbf{x}$: High probability of label 1
- High negative score $\mathbf{w}^\top \mathbf{x}$: Low prob. of label 1 (high prob. of label 0)

Logistic Regression: Parameter Estimation

- Recall, each label y_n is binary with prob. μ_n . Assume Bernoulli likelihood:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n|x_n, \mathbf{w}) = \prod_{n=1}^N \mu_n^{y_n} (1 - \mu_n)^{1-y_n}$$

$$\text{where } \mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$$

- Negative log-likelihood

$$\text{NLL}(\mathbf{w}) = -\log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N (y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n))$$

- Plugging in $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ and chugging, we get (verify yourself)

$$\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

- To do MLE for \mathbf{w} , we'll **minimize** negative log-likelihood $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}
- Important note:** $\text{NLL}(\mathbf{w})$ is convex in \mathbf{w} , so global minima can be found

MLE Estimation for Logistic Regression

- We have $\text{NLL}(\mathbf{w}) = -\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$
- Taking the derivative of $\text{NLL}(\mathbf{w})$ w.r.t. \mathbf{w}

$$\begin{aligned}\frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[-\sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= -\sum_{n=1}^N \left(y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right)\end{aligned}$$

- Can't get a closed form estimate for \mathbf{w} by setting the derivative to zero
- One solution: Iterative minimization via gradient descent. Gradient is:

$$\mathbf{g} = \frac{\partial \text{NLL}(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

- Intuitively, a large error on $\mathbf{x}_n \Rightarrow (y_n - \mu_n)$ will be large \Rightarrow large contribution (positive/negative) of \mathbf{x}_n to the gradient

MLE Estimation via Gradient Descent

- Gradient descent (GD) or steepest descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t$$

where η_t is the learning rate (or step size), and \mathbf{g}_t is gradient at step t

- GD can converge slowly and is also sensitive to the step size
- Several ways to remedy this¹. E.g.,
 - Choose the optimal step size η_t by **line-search**
 - Add a **momentum term** to the updates

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t + \alpha_t (\mathbf{w}_t - \mathbf{w}_{t-1})$$

- Use methods such as **conjugate gradient**
- Use **second-order methods** (e.g., **Newton's method**) to exploit the curvature of the objective function $\text{NLL}(\mathbf{w})$: Require the **Hessian matrix**

¹Also see: "A comparison of numerical optimizers for logistic regression" by Tom Minka

MLE Estimation via Newton's Method

- Update via Newton's method:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{H}_t^{-1} \mathbf{g}_t$$

where \mathbf{H}_t is the Hessian matrix at step t

- Hessian: double derivative of the objective function ($\text{NLL}(\mathbf{w})$ in this case)

$$\mathbf{H} = \frac{\partial^2 \text{NLL}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \frac{\partial \mathbf{g}^\top}{\partial \mathbf{w}}$$

- Recall that the gradient is: $\mathbf{g} = -\sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$
- Thus $\mathbf{H} = \frac{\partial \mathbf{g}^\top}{\partial \mathbf{w}} = -\frac{\partial}{\partial \mathbf{w}} \sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n^\top = \sum_{n=1}^N \frac{\partial \mu_n}{\partial \mathbf{w}} \mathbf{x}_n^\top$
- Using the fact that $\frac{\partial \mu_n}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \left(\frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)} \right) = \mu_n(1 - \mu_n) \mathbf{x}_n$, we have

$$\mathbf{H} = \sum_{n=1}^N \mu_n(1 - \mu_n) \mathbf{x}_n \mathbf{x}_n^\top = \mathbf{X}^\top \mathbf{S} \mathbf{X}$$

where \mathbf{S} is a diagonal matrix with its n^{th} diagonal element $= \mu_n(1 - \mu_n)$

MLE Estimation via Newton's Method

- Update via Newton's method:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \mathbf{H}_t^{-1} \mathbf{g}_t \\ &= \mathbf{w}_t - (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^\top (\boldsymbol{\mu}_t - \mathbf{y}) \\ &= \mathbf{w}_t + (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}_t) \\ &= (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} [(\mathbf{X}^\top \mathbf{S}_t \mathbf{X}) \mathbf{w}_t + \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}_t)] \\ &= (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^\top [\mathbf{S}_t \mathbf{X} \mathbf{w}_t + \mathbf{y} - \boldsymbol{\mu}_t] \\ &= (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{S}_t [\mathbf{X} \mathbf{w}_t + \mathbf{S}^{-1} (\mathbf{y} - \boldsymbol{\mu}_t)] \\ &= (\mathbf{X}^\top \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{S}_t \hat{\mathbf{y}}_t \end{aligned}$$

- Interpreting the solution found by Newton's method:
 - It basically solves an **Iteratively Reweighted Least Squares (IRLS)** problem

$$\arg \min_{\mathbf{w}} \sum_{n=1}^N S_{tn} (\hat{\mathbf{y}}_{tn} - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Note that the (redefined) response vector $\hat{\mathbf{y}}_t$ changes in each iteration
- Each term in the objective has weight S_{tn} (changes in each iteration)
- The weight S_{tn} is the n^{th} diagonal element of \mathbf{S}_t

MAP Estimation for Logistic Regression

- MLE estimate of \mathbf{w} can lead to overfitting. Solution: use a prior on \mathbf{w}
- Just like the linear regression case, let's put a Gaussian prior on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- MAP objective: MLE objective + $\log p(\mathbf{w})$
- Leads to the objective (negative of log posterior, ignoring constants):

$$\text{NLL}(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Estimation of \mathbf{w} proceeds the same way as MLE except that now we have

$$\text{Gradient: } \mathbf{g} = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}) + \lambda \mathbf{w}$$

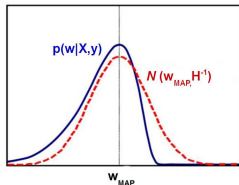
$$\text{Hessian: } \mathbf{H} = \mathbf{X}^\top \mathbf{S} \mathbf{X} + \lambda \mathbf{I}_D$$

- Can now apply iterative optimization (gradient des., Newton's method, etc.)
- **Note:** MAP estimation for log. reg. is equivalent to **regularized log. reg.**

Fully Bayesian Estimation for Logistic Regression

- What about the **full posterior** on \mathbf{w} ?
- Not as easy to estimate as in the linear regression case!
- Reason: likelihood (logistic-Bernoulli) and prior (Gaussian) not conjugate
- Need to *approximate* the posterior in this case
- A crude approximation: **Laplace approximation**: Approximate a posterior by a **Gaussian** with **mean = MAP estimate** and **covariance = inverse hessian**

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



- Will see other ways of approximating the posterior later during the semester

Derivation of the Laplace Approximation

- The posterior $p(\mathbf{w}|\mathbf{X}, y) = \frac{p(y|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(y|\mathbf{X})}$. Let's approximate it as

$$p(\mathbf{w}|\mathbf{X}, y) = \frac{\exp(-E(\mathbf{w}))}{Z}$$

where $E(\mathbf{w}) = -\log p(y|\mathbf{X}, \mathbf{w})p(\mathbf{w})$ and Z is the normalizer

- Expand $E(\mathbf{w})$ around its minima ($\mathbf{w}_* = \mathbf{w}_{MAP}$) using 2^{nd} order Taylor exp.

$$\begin{aligned} E(\mathbf{w}) &\approx E(\mathbf{w}_*) + (\mathbf{w} - \mathbf{w}_*)^\top \mathbf{g} + \frac{1}{2}(\mathbf{w} - \mathbf{w}_*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}_*) \\ &= E(\mathbf{w}_*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}_*) \quad (\text{because } \mathbf{g} = 0 \text{ at } \mathbf{w}_*) \end{aligned}$$

- Thus the posterior

$$p(\mathbf{w}|\mathbf{X}, y) \approx \frac{\exp(-E(\mathbf{w}_*)) \exp(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}_*))}{Z}$$

- Using $\int_{\mathbf{w}} p(\mathbf{w}|\mathbf{X}, y) d\mathbf{w} = 1$, we get $Z = \exp(-E(\mathbf{w}_*))(2\pi)^{D/2} |\mathbf{H}|^{-1/2}$. Thus

$$p(\mathbf{w}|\mathbf{X}, y) = \mathcal{N}(\mathbf{w}_*, \mathbf{H}^{-1})$$

Multinomial Logistic Regression

- Logistic reg. can be extended to handle $K > 2$ classes)
- In this case, $y_n \in \{0, 1, 2, \dots, K - 1\}$ and label probabilities are defined as

$$p(y_n = k | \mathbf{x}_n, \mathbf{W}) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{\ell=1}^K \exp(\mathbf{w}_\ell^\top \mathbf{x}_n)} = \mu_{nk}$$

- μ_{nk} : probability that example n belongs to class k . Also, $\sum_{\ell=1}^K \mu_{n\ell} = 1$
- $\mathbf{W} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_K]$ is $D \times K$ **weight matrix** (column k for class k)
- Likelihood for the **multinomial (or multinoulli) logistic regression** model

$$p(\mathbf{y} | \mathbf{X}, \mathbf{W}) = \prod_{n=1}^N \prod_{\ell=1}^K \mu_{n\ell}^{y_{n\ell}}$$

where $y_{n\ell} = 1$ if true class of example n is ℓ and $y_{n\ell'} = 0$ for all other $\ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for \mathbf{W} similar to the binary case
- **Decision rule:** $y_* = \arg \max_{\ell=1, \dots, K} \mathbf{w}_\ell^\top \mathbf{x}_*$, i.e., predict the class whose weight vector gives the largest score (or, equivalently, the largest probability)

Next class: Generalized Linear Models