

Machine Learning - Waseda University

Logistic Regression

AD

June 2011

- Assume you are given some training data $\{\mathbf{x}^i, y^i\}_{i=1}^N$ where $\mathbf{x}^i \in \mathbb{R}^d$ and y^i can take C different values.
- Given an input test data \mathbf{x} , you want to predict/estimate the output y associated to \mathbf{x} .
- A common approach consists of using

$$p(y = k | \mathbf{x}) = \frac{p(\mathbf{x} | y = k) p(y = k)}{\sum_{j=0}^{C-1} p(\mathbf{x} | y = j) p(y = j)}.$$

- This requires modelling and learning the parameters of the class conditional density of features $p(\mathbf{x} | y = k)$.
- This can be tedious for complicated problems.

Logistic Regression

- Discriminative model: we model and learn directly $p(y = k | \mathbf{x})$ and bypassing the introduction of $p(\mathbf{x} | y = k)$.
- Consider the following model for $C = 2$ (binary classification)

$$\begin{aligned} p(y = 1 | \mathbf{x}, \mathbf{w}) &= 1 - p(y = 0 | \mathbf{x}, \mathbf{w}) \\ &= g(\mathbf{w}^T \mathbf{x}) \end{aligned}$$

where $\mathbf{w} = (w_0 \cdots w_d)^T$, $\mathbf{x} = (x_0 \cdots x_d)^T$ so

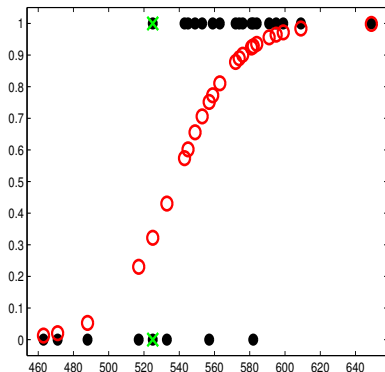
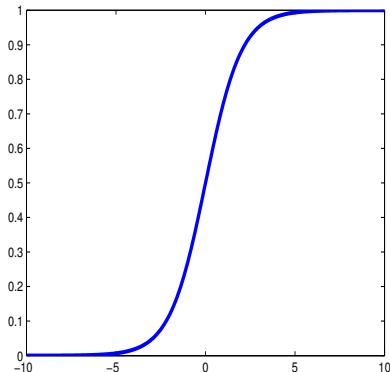
$$z = \mathbf{w}^T \mathbf{x} = w_0 + \sum_{j=1}^d w_j x_j$$

and g is a “squashing” function: $g : \mathbb{R} \rightarrow [0, 1]$.

- Logistic regression corresponds to

$$g(z) = \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{1 + \exp(z)}.$$

Logistic Function



(Left) logistic or sigmoid function (Right) logistic regression for x =SAT score and y =pass/fail class (solid black dots are the data), open red circles are predicted probabilities.

Logistic Regression

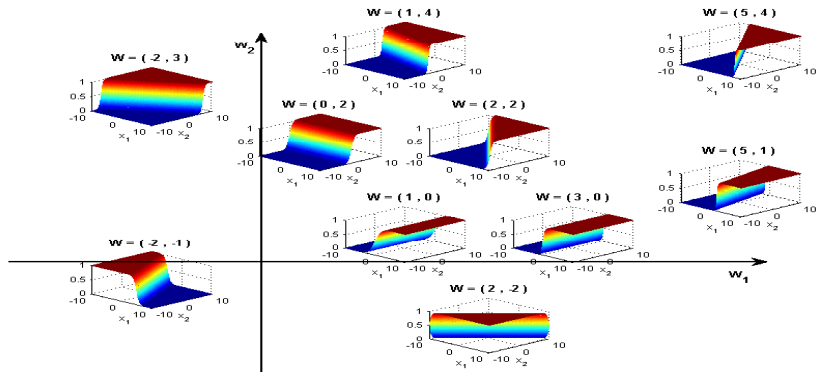
- The log odds ratio satisfies

$$LOR(\mathbf{x}) = \log \frac{p(y = 1 | \mathbf{x}, \mathbf{w})}{p(y = 0 | \mathbf{x}, \mathbf{w})} = \mathbf{w}^T \mathbf{x}$$

so the logistic parameters are easily interpretable.

- If $w_j > 0$, then increasing x_j makes $y = 1$ more likely while decreasing x_j makes $y = 0$ more likely (and opposite if $w_j = 0$). $w_j = 0$ means x_j has no impact on the outcome.
- Logistic regression partitions the input space into two regions whose decision boundary is $\{\mathbf{x} : LOR(\mathbf{x}) = 0\} = \{\mathbf{x} : \mathbf{w}^T \mathbf{x} = 0\}$
- Simple model of a neuron: it forms a weighted sum of its inputs and the “fires” an output pulse if this sum exceeds a threshold. Logistic regression mimics this as you can sort of think of it as a process which “fires” if $p(y = 1 | \mathbf{x}, \mathbf{w}) > p(y = 0 | \mathbf{x}, \mathbf{w})$ equivalently if $LOR(\mathbf{x}) > 0$.

Logistic Function in Two Dimensions



Plots of $p(y=1 | w_1x_1 + w_2x_2)$. Here $\mathbf{w} = (w_1, w_2)$ define the normal to the decision boundary. Points to the right have $\mathbf{w}^T \mathbf{x} > 0$ and to the left have $\mathbf{w}^T \mathbf{x} < 0$.

Using Basis Functions for Logistic Regression

- Similarly to regression, we can use basis functions; i.e.

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = g(\mathbf{w}^T \Phi(\mathbf{x}))$$

where $\mathbf{w} = (w_1 \ \cdots \ w_m)^T$, $\Phi(\mathbf{x}) = (\Phi_1(\mathbf{x}) \ \cdots \ \Phi_m(\mathbf{x}))^T$.

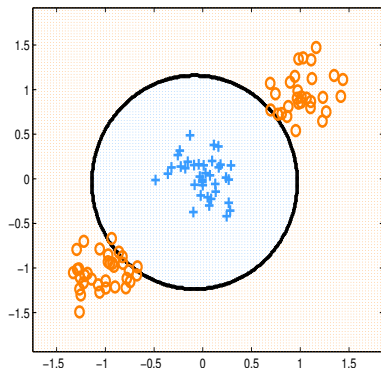
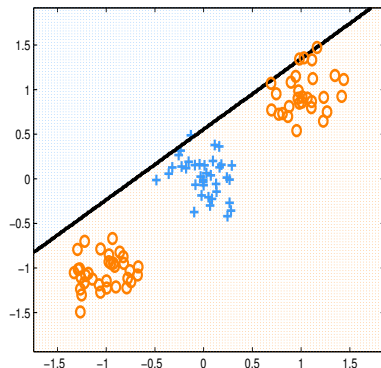
- For example, if $\mathbf{x} \in \mathbb{R}$ then we can pick

$$\Phi(\mathbf{x}) = (1, x, \dots, x^m)$$

- For $\mathbf{x} \in \mathbb{R}^d$, we can pick some radial basis functions

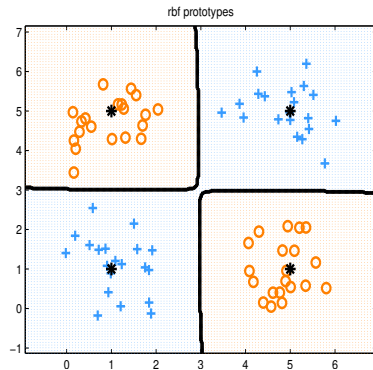
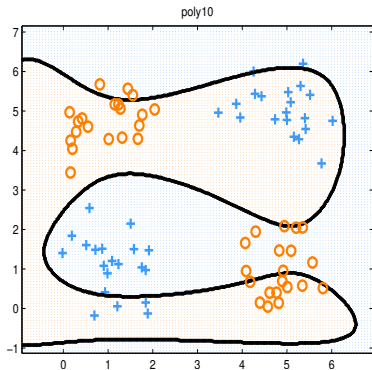
$$\Phi_j(\mathbf{x}) = \exp\left(-\frac{(\mathbf{x} - \mu_j)^T (\mathbf{x} - \mu_j)}{2\sigma^2}\right).$$

Example



(left) Logistic regression in the original feature space $\mathbf{x} = (x_1, x_2)$. (right) Logistic regression obtained after performing a 2nd degree poly expansion $\Phi(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_2^2)$.

Example



(left) Logistic regression for $\Phi(\mathbf{x}) = (1, x_1, x_2, \dots, x_1^{10}, x_2^{10})$. (right) Logistic regression using 4 radial basis functions with centers μ_j specified by black crosses.

MLE Parameter Learning for Logistic Regression

- To learn the parameters \mathbf{w} , we can maximize w.r.t \mathbf{w} the (conditional) log-likelihood function

$$\begin{aligned}l(\mathbf{w}) &= \log p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \mathbf{w}\right) = \log \prod_{i=1}^N p(y^i \mid \mathbf{x}^i, \mathbf{w}) \\&= \sum_{i=1}^N \log p(y^i \mid \mathbf{x}^i, \mathbf{w})\end{aligned}$$

- We have

$$\begin{aligned}l(\mathbf{w}) &= \sum_{i=1}^N y^i \log p(y^i = 1 \mid \mathbf{x}^i, \mathbf{w}) + (1 - y^i) \log p(y^i = 0 \mid \mathbf{x}^i, \mathbf{w}) \\&= - \sum_{i=1}^N (1 - y^i) \mathbf{w}^T \Phi(\mathbf{x}^i) - \sum_{i=1}^N \log \left(1 + \exp\left(-\mathbf{w}^T \Phi(\mathbf{x}^i)\right)\right)\end{aligned}$$

- **Good news:** $l(\mathbf{w})$ is concave so there is no local maxima.
- **Bad news:** there is no-closed form solution for $\hat{\mathbf{w}}_{\text{MLE}}$.

- Gradient ascent is one of the most basic method to maximize a function.
- It is an iterative procedure such that at iteration t :

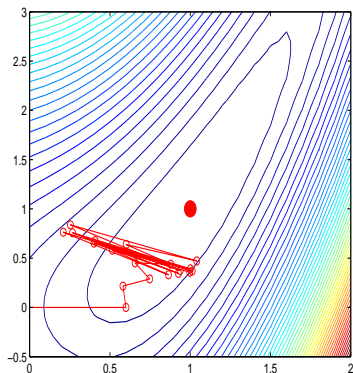
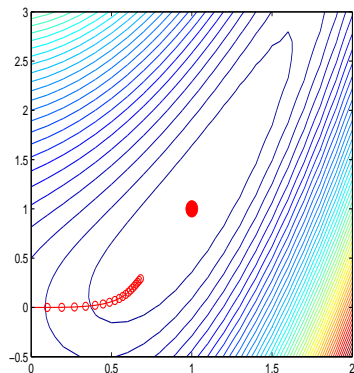
$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + \eta \nabla_{\mathbf{w}} I(\mathbf{w})|_{\mathbf{w}^{(t-1)}}$$

where the gradient is

$$\nabla_{\mathbf{w}} I(\mathbf{w}) = \left[\frac{\partial I(\mathbf{w})}{\partial w_1} \quad \dots \quad \frac{\partial I(\mathbf{w})}{\partial w_m} \right]^T$$

and $\eta > 0$ is the learning rate.

Gradient Descent Example



Gradient descent on a simple function, starting from (0,0) for 20 steps using $\eta = 0.1$ (left) and $\eta = 0.6$ (right)

Gradient Ascent for Logistic Regression

- We have

$$\nabla l(\mathbf{w}) = \sum_{i=1}^N \left(y^i - g(\mathbf{w}^\top \Phi(\mathbf{x}^i)) \right) \Phi(\mathbf{x}^i) = \Phi^\top (\mathbf{y} - \boldsymbol{\mu})$$

where $[\Phi]_{i,j} = \Phi_j(\mathbf{x}^i)$, $\mathbf{y} = (y^1 \dots y^N)^\top$ and
 $\boldsymbol{\mu} = (g(\mathbf{w}^\top \Phi(\mathbf{x}^1)) \dots g(\mathbf{w}^\top \Phi(\mathbf{x}^N)))^\top$.

- So in vector-form, we will do

$$\begin{aligned} \mathbf{w}^{(t)} &= \mathbf{w}^{(t-1)} + \eta \nabla_{\mathbf{w}} l(\mathbf{w})|_{\mathbf{w}^{(t-1)}} \\ &= \mathbf{w}^{(t-1)} + \eta \Phi^\top (\mathbf{y} - \boldsymbol{\mu}^{(t-1)}) \end{aligned}$$

where $\boldsymbol{\mu}^{(t-1)}$ corresponds to $\boldsymbol{\mu}$ computed using $\mathbf{w}^{(t-1)}$.

Iterative Reweighted Least Squares

- Newton's method is a generic (second order) optimization algorithm which converges faster than the simple gradient algorithm. It proceeds as follows at iteration t

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \left[\nabla^2 I \left(\mathbf{w}^{(t-1)} \right) \right]^{-1} \nabla I \left(\mathbf{w}^{(t-1)} \right).$$

- We have

$$\nabla^2 I(\mathbf{w}) = -\Phi^T U \Phi$$

with U a diagonal matrix with diagonal element

$$[U]_{i,i} = g \left(\mathbf{w}^T \Phi(\mathbf{x}^i) \right) \left[1 - g \left(\mathbf{w}^T \Phi(\mathbf{x}^i) \right) \right].$$

- It can be written as

$$\begin{aligned} \mathbf{w}^{(t)} &= \left(\Phi^T U^{(t-1)} \Phi \right)^{-1} \Phi^T U^{(t-1)} \\ &\times \left(\Phi \mathbf{w}^{(t-1)} + \left[U^{(t-1)} \right]^{-1} \left(\mathbf{y} - \boldsymbol{\mu}^{(t-1)} \right) \right) \end{aligned}$$

where $U^{(t-1)}$ and $\boldsymbol{\mu}^{(t-1)}$ corresponds to U and $\boldsymbol{\mu}$ with $\mathbf{w}^{(t-1)}$.

Regularized Logistic Regression via Gaussian Prior

- Similarly to regression, we can regularize the solution by assigning a Gaussian prior to \mathbf{w}

$$p(\mathbf{w}) = \prod_{k=1}^m p(w_k) = \prod_{k=1}^m \mathcal{N}(w_k; 0, \lambda)$$

- This pushes the parameters \mathbf{w} towards zero and can prevent overfitting. In this case, we have

$$\begin{aligned}\mathbf{w}_{MAP} &= \arg \max p(\mathbf{w} | \{\mathbf{x}^i, y^i\}_{i=1}^N) \\ &= \arg \max l(\mathbf{w}) - \frac{\mathbf{w}^T \mathbf{w}}{2\lambda}.\end{aligned}$$

- \mathbf{w}_{MAP} can be computed iteratively using

$$\mathbf{w}^{(t-1)} + \eta \left\{ -\lambda^{-1} \mathbf{w}^{(t-1)} + \Phi^T (\mathbf{y} - \boldsymbol{\mu}^{(t-1)}) \right\}$$

- Regularization parameter can be estimated using cross-validation or by maximizing marginal likelihood.

Regularized Logistic Regression via Laplace Prior

- Similarly to regression, we can regularize the solution by assigning a Gaussian prior to \mathbf{w}

$$p(\mathbf{w}) = \prod_{k=1}^m p(w_k) = \prod_{k=1}^m \frac{1}{2\lambda} \exp\left(-\frac{|w_k|}{\lambda}\right)$$

- This pushes the parameters \mathbf{w} towards zero and can prevent overfitting. In this case, we have

$$\begin{aligned}\mathbf{w}_{MAP} &= \arg \max_{\mathbf{w}} p(\mathbf{w} | \{\mathbf{x}^i, y^i\}_{i=1}^N) \\ &= \arg \max_{\mathbf{w}} l(\mathbf{w}) - \frac{1}{2\lambda} \sum_{k=1}^m |w_k|.\end{aligned}$$

- The objective function is convex and efficient procedures have been developed to compute \mathbf{w}_{MAP} . Similarly to the regression case, this can lead to sparse solution; e.g. you can have $w_{k,MAP} = 0$ exactly.
- Regularization parameter can be estimated using cross-validation or by maximizing marginal likelihood.

Multinomial Logistic Regression

- Consider now the case where $C > 2$. We could consider the following generalization

$$p\left(y = c \mid \mathbf{x}, \{\mathbf{w}_c\}_{c=1}^C\right) = \frac{\exp\left(\mathbf{w}_c^T \Phi(\mathbf{x})\right)}{\sum_{k=1}^C \exp\left(\mathbf{w}_k^T \Phi(\mathbf{x})\right)} \text{ for } c = 1, \dots, C$$

but this is not identifiable as

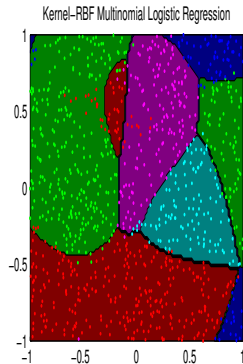
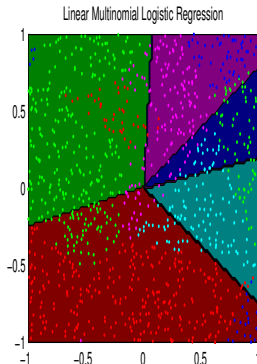
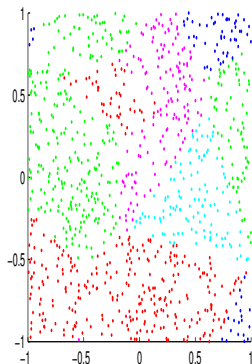
$$p\left(y = c \mid \mathbf{x}, \{\mathbf{w}_c + \mathbf{w}'\}_{c=1}^C\right) = p\left(y = c \mid \mathbf{x}, \{\mathbf{w}_c\}_{c=1}^C\right).$$

- Hence we set $\mathbf{w}_C = (0 \ \cdots \ 0)^T$ to obtain

$$p\left(y = c \mid \mathbf{x}, \{\mathbf{w}_c\}_{c=1}^{C-1}\right) = \frac{\exp\left(\mathbf{w}_c^T \Phi(\mathbf{x})\right)}{1 + \sum_{k=1}^{C-1} \exp\left(\mathbf{w}_k^T \Phi(\mathbf{x})\right)} \text{ for } c = 1, \dots, C-1$$
$$p\left(y = C \mid \mathbf{x}, \{\mathbf{w}_c\}_{c=1}^{C-1}\right) = \frac{1}{1 + \sum_{k=1}^{C-1} \exp\left(\mathbf{w}_k^T \Phi(\mathbf{x})\right)}.$$

- The (conditional) log-likelihood is concave w.r.t $\{\mathbf{w}_c\}_{c=1}^{C-1}$ so MLE estimates can be computed using gradient.

Example



(left) Some 5 class data in 2d (center) Multinomial logistic regression in the original feature space $\mathbf{x} = (x_1, x_2)$ (right) RBF basis functions with bandwidth 1 using $m = 1 + N$. We use all the data points as centers.

Full Bayesian Analysis of Logistic Regression

- Even for Gaussian priors on \mathbf{w} , one cannot compute

$$p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \lambda\right) = \frac{p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \mathbf{w}\right) p(\mathbf{w} \mid \lambda)}{p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \lambda\right)}$$

where

$$p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \lambda\right) = \int p\left(\{y^i\}_{i=1}^N \mid \{\mathbf{x}^i\}_{i=1}^N, \mathbf{w}\right) p(\mathbf{w} \mid \lambda) d\mathbf{w}$$

- Contrary to regression, there is no closed form Bayesian analysis possible.
- If you want to do Bayesian inference, then approximations are necessary.