

Probabilistic Models for Supervised Learning(2): Logistic and Softmax Regression

CS771: Introduction to Machine Learning

Piyush Rai

Logistic Regression (LR)

Multi-class extension known as “softmax regression”

Both very widely used

The word “regression” is a misnomer. Both are classification models

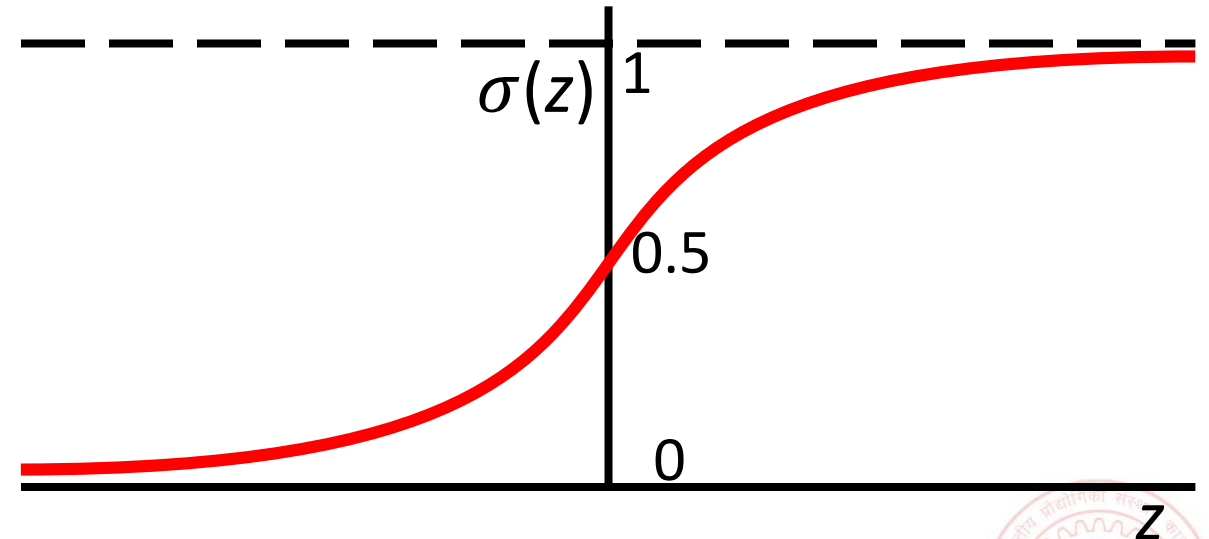
2



- A probabilistic model for binary classification
- Learns the PMF of the output label given the input, i.e., $p(y|\mathbf{x})$
- A **discriminative model**: Does not model inputs \mathbf{x} (only relationship b/w \mathbf{x} and \mathbf{y})
- Uses the **sigmoid function** to define the conditional probability of \mathbf{y} being 1

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

A linear model



- Here $\mathbf{w}^\top \mathbf{x}$ is the score for input \mathbf{x} . The sigmoid turns it into a probability

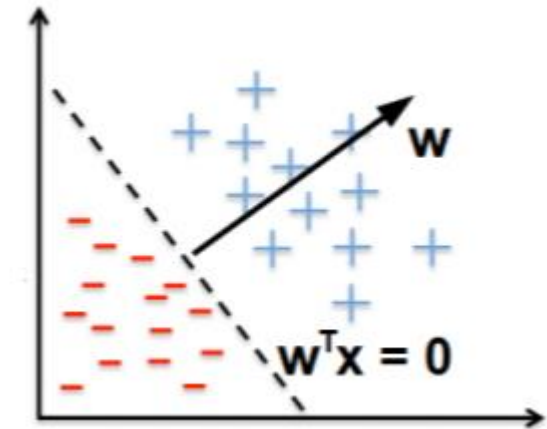


LR: Decision Boundary

- At the decision boundary where both classes are equiprobable

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

A linear hyperplane



- Very large positive $\mathbf{w}^\top \mathbf{x}$ means $p(y = 1|\mathbf{w}, \mathbf{x})$ close to 1
- Very large negative $\mathbf{w}^\top \mathbf{x}$ means $p(y = 0|\mathbf{w}, \mathbf{x})$ close to 1
- At decision boundary, $\mathbf{w}^\top \mathbf{x} = 0$ implies $p(y = 1|\mathbf{w}, \mathbf{x}) = p(y = 0|\mathbf{w}, \mathbf{x}) = 0.5$



MLE for Logistic Regression

Assumed 0/1, not -1/+1

- Likelihood (PMF of each input's label) is Bernoulli with prob $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$
 $p(y_n | \mathbf{w}, \mathbf{x}_n) = \text{Bernoulli}(\mu_n) = \mu_n^{y_n} (1 - \mu_n)^{1 - y_n}$
- Overall likelihood, assuming i.i.d. observations

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

- The negative log-likelihood $NLL(\mathbf{w}) = -\log p(\mathbf{y} | \mathbf{w}, \mathbf{X})$ simplifies to

"cross-entropy" loss (a popular loss function for classification)

Loss function

$$NLL(\mathbf{w}) = \sum_{n=1}^N -[y_n \log \mu_n + (1 - y_n) \log (1 - \mu_n)]$$

Very large loss if y_n close to 1 and μ_n close to 0, or vice-versa

- Plugging in $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$ and simplifying

Good news: For LR, NLL is convex

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

No closed-form expression for $\hat{\mathbf{w}}_{MLE} = \arg \min_{\mathbf{w}} NLL(\mathbf{w})$

Iterative opt needed (gradient or Hessian based). **Exercise:** Try working out the gradient of NLL and notice the expression's form



An Alternate Notation

- If we assume the label y_n as -1/+1 (not 0/1), the likelihood can be written as

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

- Slightly more convenient notation: A single expression gives the probabilities of both possible label values
- In this case, the total negative log-likelihood will be

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



MAP Estimation for Logistic Regression

- Need a prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$
- Just like probabilistic linear regression, can use a zero-mean Gaussian prior

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

Or NLL – log of prior

- The MAP objective (log of posterior) will be log-likelihood + log of prior
- Therefore the MAP solution (ignoring terms that don't depend on \mathbf{w}) will be

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

Good news:
convex objective

- Just like MLE case, no closed form solution. Iterative opt methods needed
 - Highly efficient solvers (both first and second order) exist for MLE/MAP estimation for LR



Fully Bayesian Inference for Logistic Regression

- Doing fully Bayesian inference would require computing the posterior

Gaussian

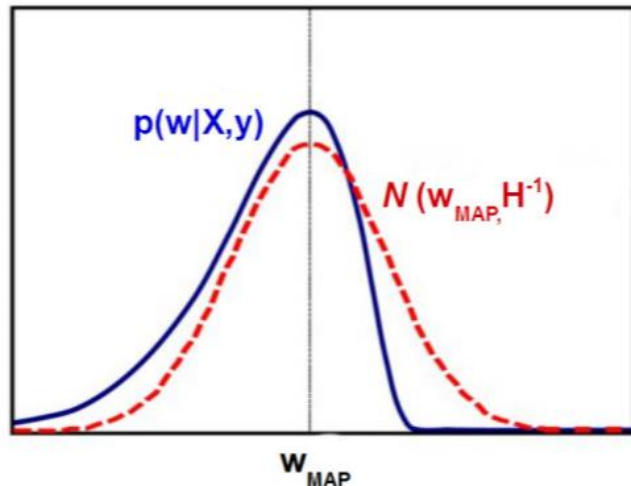
Bernoulli

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{w})p(\mathbf{y}|\mathbf{X}, \mathbf{w})}{p(\mathbf{y}|\mathbf{X})} = \frac{p(\mathbf{w}) \prod_{n=1}^N p(y_n|\mathbf{w}, \mathbf{x}_n)}{\int p(\mathbf{w}) \prod_{n=1}^N p(y_n|\mathbf{w}, \mathbf{x}_n) d\mathbf{w}}$$

Unfortunately, Gaussian and Bernoulli are not conjugate with each other, so analytic expression for the posterior can't be obtained unlike prob. linear regression



- Need to approximate the posterior in this case
- We will use a simple approximation called Laplace approximation



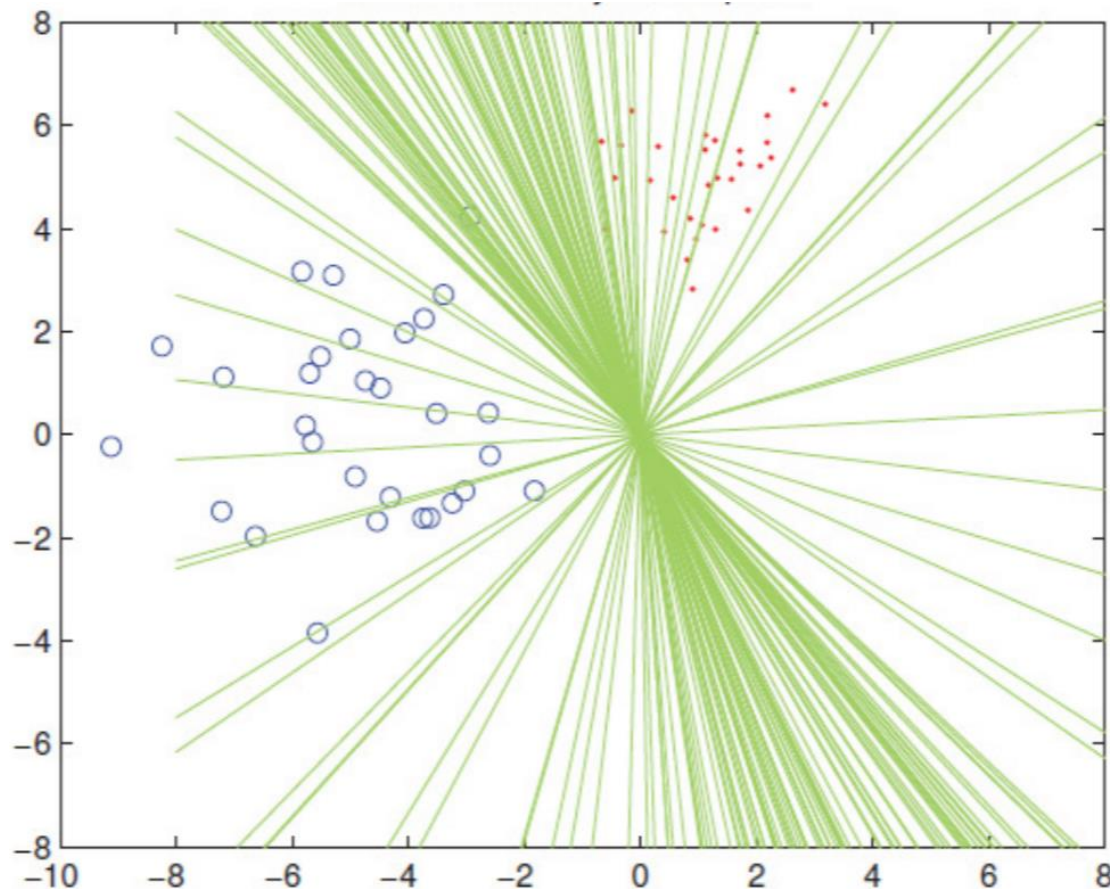
Approximates the posterior of \mathbf{w} by a Gaussian whose mean is the MAP solution $\hat{\mathbf{w}}_{MAP}$ and covariance matrix is the inverse of the Hessian (Hessian: **second derivative of the negative log-posterior of the LR model**)

Can also employ more advanced posterior approximation methods, like MCMC and variational inference (beyond the scope of CS771)



Posterior for LR: An Illustration

- Can sample from the posterior of the LR model
- Each sample will give a weight vec defining a hyperplane separator



Not all separators are equally good; their goodness depends on their posterior probabilities

When making predictions, we can still use all of them but weighted by their importance based on their posterior probabilities

That's exactly what we do when computing the predictive distribution



Logistic Regression: Predictive Distribution

- When using MLE/MAP solution $\hat{\mathbf{w}}_{opt}$, can use the **plug-in predictive distribution**

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ \approx p(y_* = 1 | \hat{\mathbf{w}}_{opt}, \mathbf{x}_*) = \sigma(\hat{\mathbf{w}}_{opt}^\top \mathbf{x}_*)$$

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \text{Bernoulli}[\sigma(\hat{\mathbf{w}}_{opt}^\top \mathbf{x}_*)]$$

- When using fully Bayesian inference, we must compute the posterior predictive

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Integral not tractable and must be approximated

sigmoid

Gaussian (if using Laplace approx.)

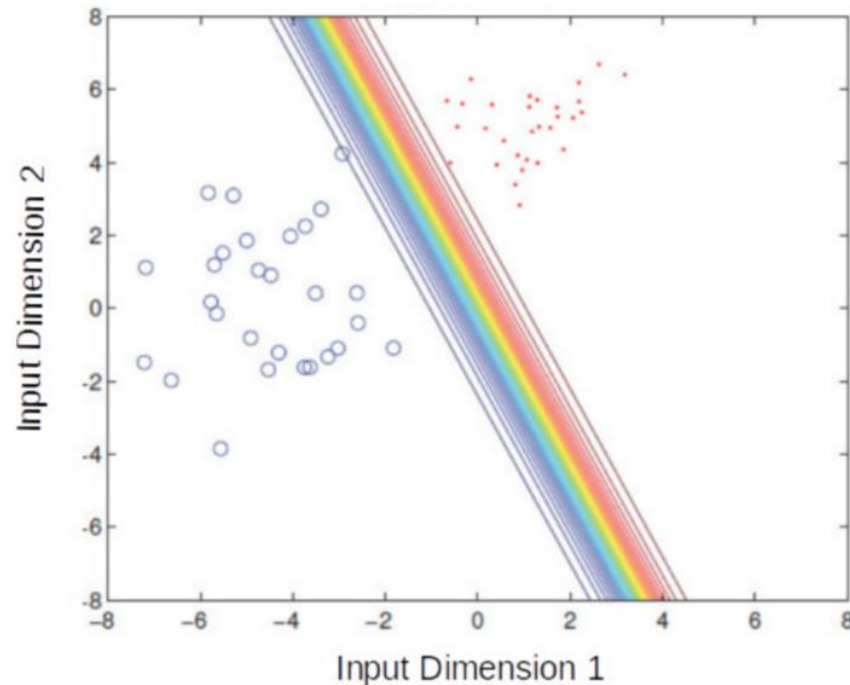
Monte-Carlo approximation of this integral is one possible way

Generate M samples $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M$, from the Gaussian approx. of posterior and use $p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M p(y_* = 1 | \mathbf{w}_m, \mathbf{x}_*) = \frac{1}{M} \sum_{m=1}^M \sigma(\mathbf{w}_m^\top \mathbf{x}_*)$

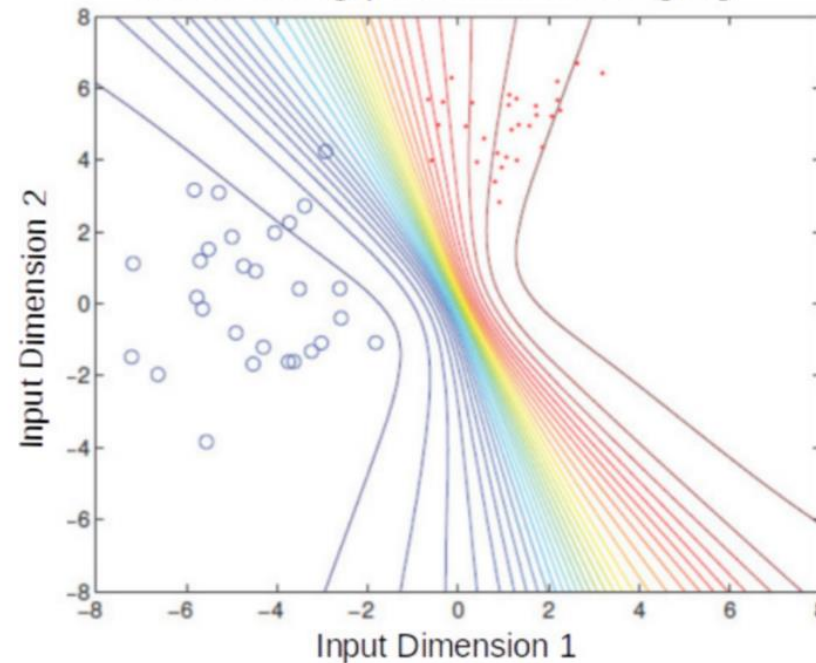
LR: Plug-in Prediction vs Posterior Averaging

10

Logistic Regression decision boundary
when using a point estimate of w



Logistic Regression decision boundary
when using posterior averaging



Posterior averaging is like using an ensemble of models. In this example, each model is a linear classifier but the ensemble-like effect resulted in nonlinear boundaries



Multiclass Logistic (a.k.a. Softmax) Regression

- Also called **multinoulli/multinomial regression**: Basically, LR for $K > 2$ classes
- In this case, $y_n \in \{1, 2, \dots, K\}$ and label probabilities are defined as

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

Softmax function

Also note that $\sum_{\ell=1}^K \mu_{n\ell} = 1$ for any input \mathbf{x}_n

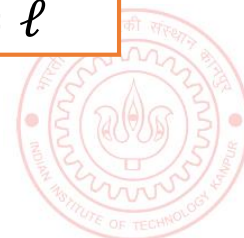


- K weight vecs $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ (one per class), each D -dim, and $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$
- Each likelihood $p(y_n | \mathbf{x}_n, \mathbf{W})$ is a **multinoulli** distribution. Therefore total likelihood

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

Notation: $y_{n\ell} = 1$ if true class of \mathbf{x}_n is ℓ and $y_{n\ell'} = 0 \forall \ell' \neq \ell$

- Can do MLE/MAP/fully Bayesian estimation for \mathbf{W} similar to LR model



Coming up next

- Generative models for supervised learning

