

CS6375: Machine Learning

Gautam Kunapuli

Logistic Regression



THE UNIVERSITY OF TEXAS AT DALLAS

Erik Jonsson School of Engineering and Computer Science

Naïve Bayes Classifier for Continuous Features

Naïve Bayes assumption: features are conditionally independent given y

$$P(x_1, x_2, \dots, x_d | y) = \prod_{j=1}^d P(x_j | y)$$

For a **continuous feature** x_j , we can represent $P(x_j | y = c)$ with a Gaussian distribution, with parameters μ_{jc} (mean) and σ_{jc} (standard deviation)

$$P(x_j | y = c) = \frac{1}{\sqrt{2\pi\sigma_{jc}^2}} \exp\left(-\frac{(x_j - \mu_{jc})^2}{2\sigma_{jc}^2}\right)$$

That is, if there are $c = 1, \dots, C$ classes, we model each conditional distribution with its **own** Gaussian distribution.

We can use **Maximum Likelihood Estimation** as before to estimate the parameters of the distribution:

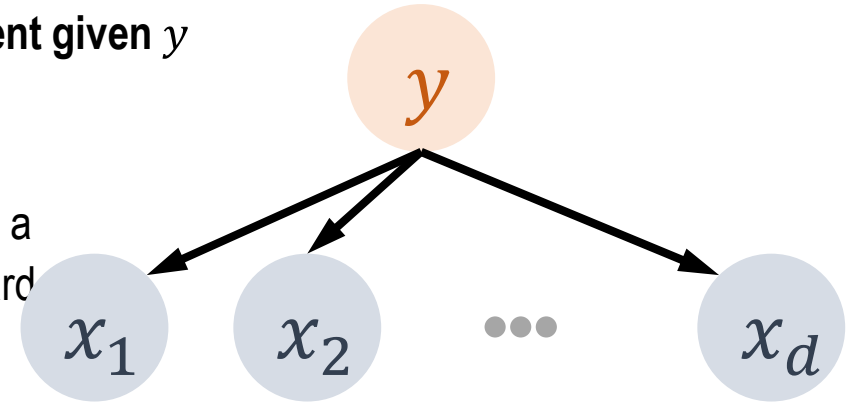
$$\mu_{jc} = \frac{1}{\sum_{i=1}^n I(y_i = c)} \sum_{i=1}^n x_{ij} \cdot I(y_i = c)$$

$$\sigma_{jc} = \frac{1}{\sum_{i=1}^n I(y_i = c)} \sum_{i=1}^n (x_{ij} - \mu_{jc})^2 \cdot I(y_i = c)$$

where i indexes training examples, j indexes features and c indexes class labels.

$I(y_i = c)$ is an indicator function that simply indicates if its argument is true or not, that is:

$$I(y_i = c) = \begin{cases} 1, & \text{if } y_i = c \\ 0, & \text{otherwise} \end{cases}$$



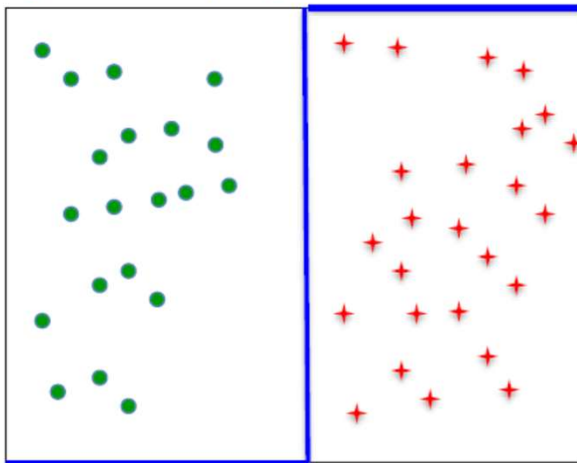
Logistic Regression

Misclassification Minimization:

Learn $p(y|\mathbf{x})$ directly from the data

- Assume a functional form, (e.g., a linear classifier $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$) such that
 - $p(y = 1 | \mathbf{x}) = 1$ on one side and
 - $p(y = 1 | \mathbf{x}) = 0$ on the other side
 that is $p(y = -1 | \mathbf{x}) = 1 - p(y = 1 | \mathbf{x}) = 1$
- Not differentiable
- Makes it difficult to learn
- Can't handle noisy labels

$$p(Y = 1 | \mathbf{x}) = 0$$

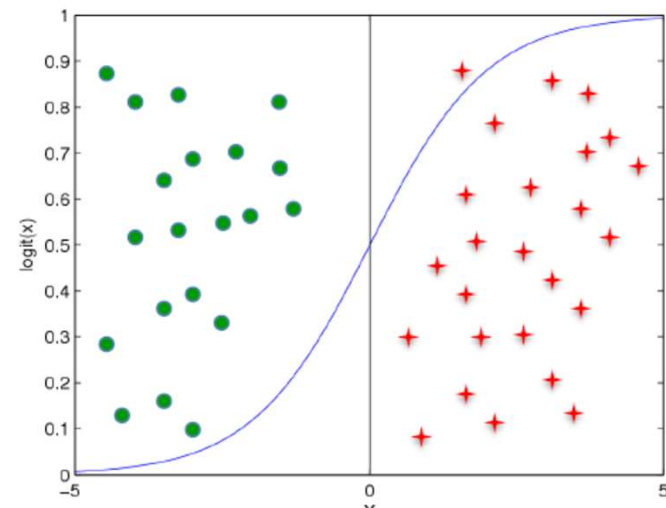


$$p(Y = 1 | \mathbf{x}) = 1$$

Logistic Loss Function:

Learn or $p(y|\mathbf{x})$ directly from the data

- Assume a functional form, (e.g., a linear classifier $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$) such that
 - $p(y = -1 | \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x} + b)}$ on one side and
 - $p(y = 1 | \mathbf{x}) = \frac{\exp(\mathbf{w}^T \mathbf{x} + b)}{1 + \exp(\mathbf{w}^T \mathbf{x} + b)}$ on the other side
 that is $p(y = -1 | \mathbf{x}) = 1 - p(y = 1 | \mathbf{x})$
- Differentiable
- Easy to learn
- Handles noisy labels naturally

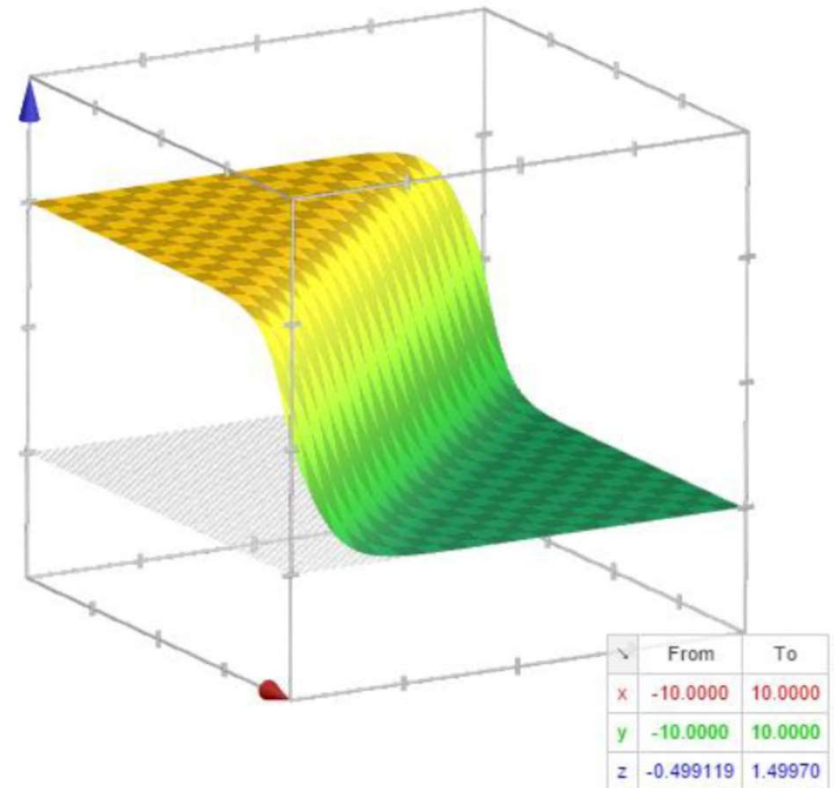
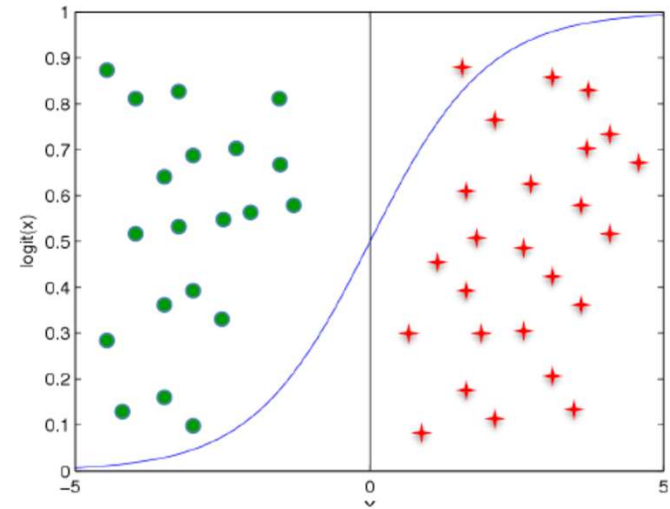


The Logistic Function

A linear function $y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ has a range from $[-\infty, \infty]$. The **logistic function** transforms this range to a probability $[0, 1]$.

Given some \mathbf{w} and b , we can classify a new point \mathbf{x} by **assigning the labels** as follows:

- $y = 1$, if $p(y = 1 | \mathbf{x}) > p(y = -1 | \mathbf{x})$ and
- $y = -1$, otherwise



The Logistic Function

A linear function $y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ has a range from $[-\infty, \infty]$. The **logistic function** transforms this range to a probability $[0, 1]$.

Given some \mathbf{w} and b , we can classify a new point \mathbf{x} by assigning the labels as follows:

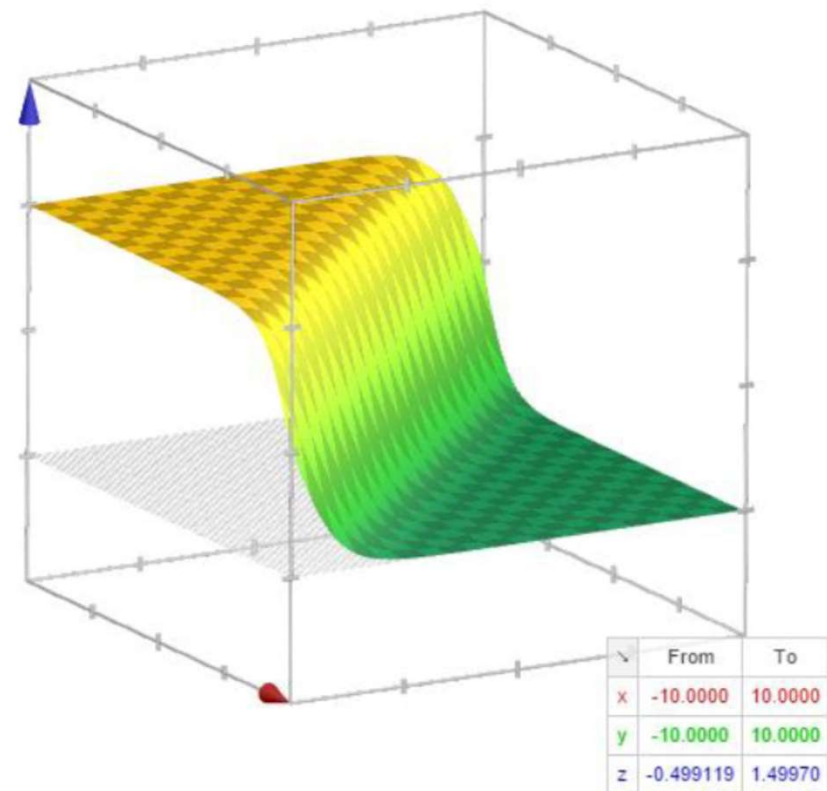
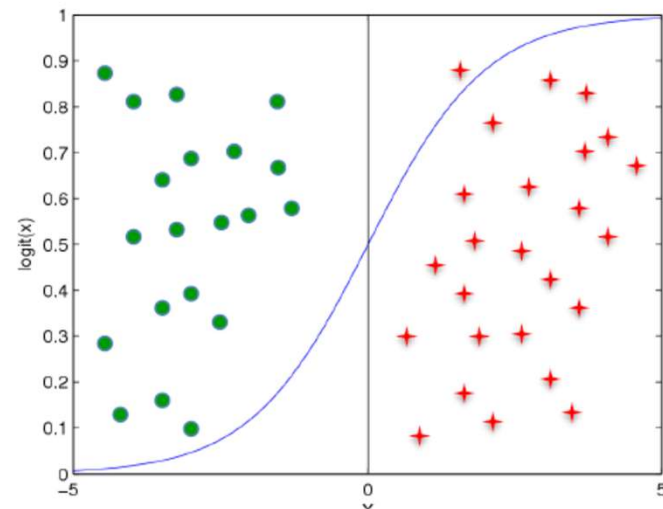
- $y = 1$, if $p(y = 1 | \mathbf{x}) > p(y = -1 | \mathbf{x})$ and
- $y = -1$, otherwise

Logistic regression **implements a linear classifier**, that is:

- $y = 1$, if $\mathbf{w}^T \mathbf{x} + b > 0$ and
- $y = -1$, otherwise

We can show this by showing the **log-odds of a training example** belonging to class $y = 1$ are:

$$\log \frac{p(y = 1 | \mathbf{x})}{p(y = -1 | \mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$



Formulating Logistic Regression

Since we are fitting a **conditional probability distribution**, we no longer minimize the loss on the training data. Instead, we are interested in **finding the distribution** h that is **most likely** given the training data.

Let $S = (\mathbf{x}_i, y_i)_{i=1}^n$ be the training data set (sample). Our goal is to find h to maximize $P(h|S)$:

$$\begin{aligned} \arg \max_h P(h | S) &= \arg \max_h \frac{P(S | h)P(h)}{P(S)} \\ &= \arg \max_h P(S | h)P(h) \\ &= \arg \max_h P(S | h) \\ &= \arg \max_h \log P(S | h) \end{aligned}$$

by **Bayes' Rule**

because $P(S)$ doesn't depend on h

assuming $P(h)$ is **uniform**

because \log is **monotonic**

In our framework, we **assume** that each training example is **identically and independently distributed (i.i.d.)**

$$\log P(S|h) = \log \prod_{i=1}^n P(\mathbf{x}_i, y_i|h) = \sum_{i=1}^n \log P(\mathbf{x}_i, y_i|h)$$

this shows that the log likelihood of a data set is the sum of the log likelihoods of the individual training examples in the data set

Learning the Weights

to make everything simpler, **assume** that the labels are $y = 1$ (for positive examples) and $y = 0$ (for negative examples, instead of $y = -1$)

$$\max_{\mathbf{w}, b} \log P(S|\mathbf{w}, b)$$

$$= \max_{\mathbf{w}, b} \log \prod_{i=1}^n P(\mathbf{x}_i, y_i|\mathbf{w}, b)$$

$$= \max_{\mathbf{w}, b} \sum_{i=1}^n \log P(\mathbf{x}_i, y_i|\mathbf{w}, b)$$

$$= \max_{\mathbf{w}, b} \sum_{i=1}^n \log [P(y_i|\mathbf{x}_i, \mathbf{w}, b) P(\mathbf{x}_i|\mathbf{w}, b)] \quad \text{by **Bayes' Rule**}$$

$$= \max_{\mathbf{w}, b} \sum_{i=1}^n \log P(y_i|\mathbf{x}_i, \mathbf{w}, b) \quad \text{because } P(\mathbf{x}_i|\mathbf{w}, b) \text{ doesn't depend on } \mathbf{w} \text{ and } b$$

$$= \max_{\mathbf{w}, b} \sum_{i=1}^n y_i \log P(y_i = 1|\mathbf{x}_i, \mathbf{w}, b) + (1 - y_i) \log P(y_i = 0|\mathbf{x}_i, \mathbf{w}, b)$$

if $y_i = 1$ the log likelihood is $\log p(y = 1 | \mathbf{x})$

if $y_i = 0$ the log likelihood is $\log p(y = 0 | \mathbf{x})$

$$= \max_{\mathbf{w}, b} \sum_{i=1}^n y_i \log \frac{P(y_i = 1|\mathbf{x}_i, \mathbf{w}, b)}{P(y_i = 0|\mathbf{x}_i, \mathbf{w}, b)} + \log P(y_i = 0|\mathbf{x}_i, \mathbf{w}, b)$$

Learning the Weights

to make everything simpler, **assume** that the labels are $y = 1$ (for positive examples) and $y = 0$ (for negative examples, instead of $y = -1$)

$$\begin{aligned} & \max_{\mathbf{w}, b} \sum_{i=1}^n y_i \log \frac{P(y_i = 1 | \mathbf{x}_i, \mathbf{w}, b)}{P(y_i = 0 | \mathbf{x}_i, \mathbf{w}, b)} + \log P(y_i = 0 | \mathbf{x}_i, \mathbf{w}, b) \\ &= \max_{\mathbf{w}, b} \sum_{i=1}^n y_i (\mathbf{w}^T \mathbf{x}_i + b) - \log(1 + \exp(\mathbf{w}^T \mathbf{x}_i + b)) \\ &= \max_{\mathbf{w}, b} L(\mathbf{w}, b) \quad \text{find parameters } \mathbf{w}, b \text{ to maximize the conditional log-likelihood} \end{aligned}$$

no closed-form solution!

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \sum_{i=1}^n (y_i - p(y_i = 1 | \mathbf{x}_i, \mathbf{w})) \cdot \mathbf{x}_i \\ \nabla_b L &= \sum_{i=1}^n (y_i - p(y_i = 1 | \mathbf{x}_i, \mathbf{w})) \end{aligned}$$

gradient depends $y_i - p(y_i = 1 | \mathbf{x}_i)$, the difference between the **true label** and the **predicted probability**

- if $y_i = 1$ (positive example), the gradient pushes $p(y_i = 1 | \mathbf{x}_i)$ closer to 1 (hopefully resulting in a high probability of $y_i = 1$)
- if $y_i = 0$ (negative example), the gradient pushes $p(y_i = 1 | \mathbf{x}_i)$ closer to 0 (hopefully resulting in a low probability of $y_i = 1$, which is a high probability of $y_i = 0$)

Learning the Weights

to make everything simpler, **assume** that the labels are $y = 1$ (for positive examples) and $y = 0$ (for negative examples, instead of $y = -1$)

$$\begin{aligned} & \max_{\mathbf{w}, b} \sum_{i=1}^n y_i \log \frac{P(y_i = 1 | \mathbf{x}_i, \mathbf{w}, b)}{P(y_i = 0 | \mathbf{x}_i, \mathbf{w}, b)} + \log P(y_i = 0 | \mathbf{x}_i, \mathbf{w}, b) \\ &= \max_{\mathbf{w}, b} \sum_{i=1}^n y_i (\mathbf{w}^T \mathbf{x}_i + b) - \log(1 + \exp(\mathbf{w}^T \mathbf{x}_i + b)) \\ &= \max_{\mathbf{w}, b} L(\mathbf{w}, b) \quad \text{find parameters } \mathbf{w}, b \text{ to maximize the conditional log-likelihood} \end{aligned}$$

(Batch) Gradient Ascent for Logistic Regression

Initialize: $w = w_0, b = b_0, t = 0$

Iterate until convergence

Compute updates:

$$w_{t+1} = w_t + \eta_t \nabla_{\mathbf{w}} L(f(\mathbf{x}), y)$$

$$b_{t+1} = b_t + \eta_t \nabla_b L(f(\mathbf{x}), y)$$

Check for convergence

Continue to next iteration: $t = t + 1$

- Batch algorithm: use all the data points together; $L(\mathbf{w}, b)$ is a **concave function**, so **gradient ascent** (because maximization) will converge to a **global minimum**
- Online algorithm: use (small chunks or) one data point at a time; leads to **stochastic gradient descent**, which is **highly efficient**

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \sum_{i=1}^n (y_i - p(y_i = 1 | \mathbf{x}_i, \mathbf{w})) \cdot \mathbf{x}_i \\ \nabla_b L &= \sum_{i=1}^n (y_i - p(y_i = 1 | \mathbf{x}_i, \mathbf{w})) \end{aligned}$$

gradient depends $y_i - p(y_i = 1 | \mathbf{x}_i)$, the difference between the **true label** and the **predicted probability**

- if $y_i = 1$ (positive example), the gradient pushes $p(y_i = 1 | \mathbf{x}_i)$ closer to 1 (hopefully resulting in a high probability of $y_i = 1$)
- if $y_i = 0$ (negative example), the gradient pushes $p(y_i = 1 | \mathbf{x}_i)$ closer to 0 (hopefully resulting in a low probability of $y_i = 1$, which is a high probability of $y_i = 0$)

Priors and Regularization

Overfitting the training data is possible in Logistic Regression, especially when data is very high dimensional and training data is sparse; can be avoided by adding a **prior**, which leads to a penalized log-likelihood

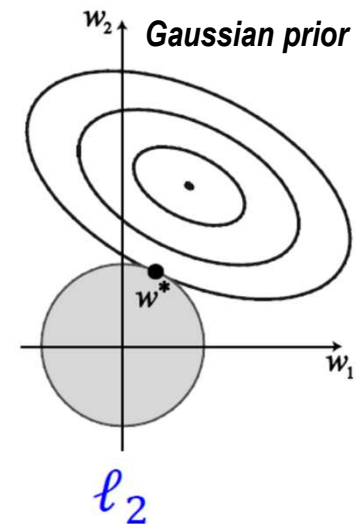
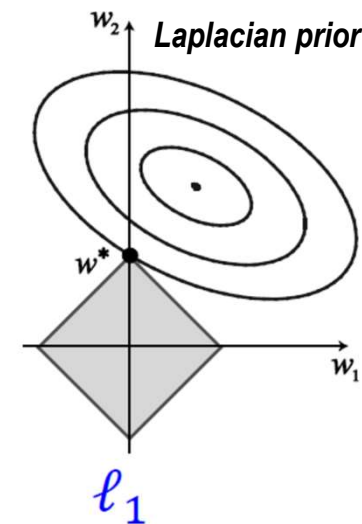
$$\max_{\mathbf{w}, b} \sum_{i=1}^n \log P(y_i | \mathbf{x}_i, \mathbf{w}, b) + \lambda \log P(\mathbf{w})$$

Consider a **prior distribution** on the weights to prevent overfitting

- assume weights from a normal distribution with zero mean, identity covariance: $P(\mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\|\mathbf{w}\|^2}{2\sigma^2}\right)$
- maximizing** $P(\mathbf{w})$ pushes weights to zero, which **minimizes the complexity** of the resulting classifier; this in turn also helps avoid large weights and overfitting
- taking the logarithm gives us $\log P(\mathbf{w}) = -\|\mathbf{w}\|^2 + \text{const}$

$$\max_{\mathbf{w}, b} \sum_{i=1}^n y_i (\mathbf{w}^T \mathbf{x}_i + b) - \log(1 + \exp(\mathbf{w}^T \mathbf{x}_i + b)) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- can **also** be solved by **gradient ascent**
 - batch algorithm still has global optimal solution
 - online algorithm with one/few data points at a time will lead to a stochastic gradient descent algorithm
- regularization parameter: $\lambda > 0$**



different priors can lead to different regularization functions

Naïve Bayes vs. Logistic Regression

Non-asymptotic analysis (for Gaussian NB)

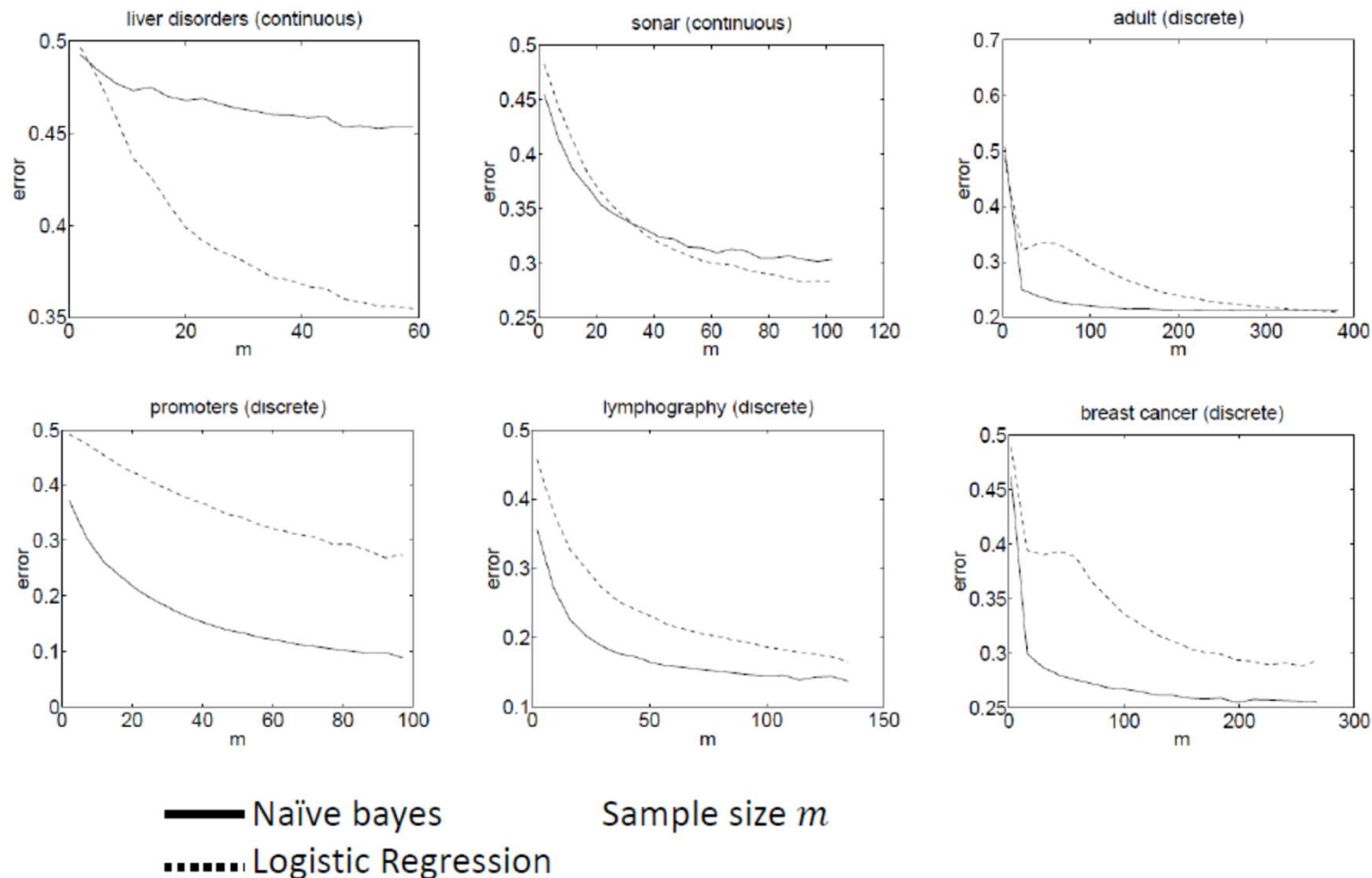
Size of training data to get close to infinite data solution, ($m = \#$ of attributes/features in X)

- **Naïve Bayes** needs $O(\log m)$ samples

*NB converges quickly to its (perhaps less helpful) asymptotic estimates; makes very strong **independence assumptions***

- **Logistic Regression** needs $O(m)$ samples

LR converges more slowly but makes no independence assumptions



Multiclass Logistic Regression

Let there be C classes: $c = 1, \dots, C$. Choose class C to be the **reference class** and represent each of the other classes as a logistic function with respect to class C :

$$\begin{array}{l}
 \log \frac{p(y = 1 | \mathbf{x})}{p(y = C | \mathbf{x})} = \mathbf{w}_1^T \mathbf{x} + b_1 \\
 \vdots \\
 \log \frac{p(y = c | \mathbf{x})}{p(y = C | \mathbf{x})} = \mathbf{w}_c^T \mathbf{x} + b_c \\
 \vdots \\
 \log \frac{p(y = C - 1 | \mathbf{x})}{p(y = C | \mathbf{x})} = \mathbf{w}_{C-1}^T \mathbf{x} + b_{C-1}
 \end{array}
 \quad
 \begin{array}{l}
 p(y = 1 | \mathbf{x}) = \frac{\exp(\mathbf{w}_1^T \mathbf{x} + b_1)}{1 + \sum_{c=1}^{C-1} \exp(\mathbf{w}_c^T \mathbf{x} + b_c)} \\
 \vdots \\
 p(y = c | \mathbf{x}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x} + b_c)}{1 + \sum_{c=1}^{C-1} \exp(\mathbf{w}_c^T \mathbf{x} + b_c)} \\
 \vdots \\
 p(y = C - 1 | \mathbf{x}) = \frac{\exp(\mathbf{w}_{C-1}^T \mathbf{x} + b_{C-1})}{1 + \sum_{c=1}^{C-1} \exp(\mathbf{w}_c^T \mathbf{x} + b_c)} \\
 \\
 p(y = C | \mathbf{x}) = \frac{1}{1 + \sum_{c=1}^{C-1} \exp(\mathbf{w}_c^T \mathbf{x} + b_c)}
 \end{array}$$

Gradient ascent can be applied **simultaneously** to train all the weight vectors and bias constants.