

LINEAR MODELS FOR CLASSIFICATION

Chia-Ping Chen

Professor

National Sun Yat-sen University

Department of Computer Science and Engineering

Machine Learning

- Discriminant Functions
- Generative Models
- Discriminative Models
- Laplace Approximation
- Bayesian Logistic Regression

Definition. The goal of **classification** is to take an input instance x and assign it to one discrete class. The assignment is achieved by a **classification function**.

- In **AlphaGo**, decide the next move given previous moves
- Given a **speech** waveform, decide its emotional category
- Given an **image** of written digit, decide the digit
- Decide whether to buy, sell, or keep for a **stock portfolio**

Definition. By a classification function, an **input space** or **feature space** is divided into **decision regions** for the classes. The boundary between adjacent decision regions is a **decision boundary** or **decision surface**.

Definition. A **linear model for classification** means that the decision boundaries are linear functions of the input vector x . That is, the decision boundaries are hyperplanes. A **generalized linear model** means that the decision boundaries are linear functions of a feature vector ϕ .

Definition. A data set whose data points of different classes can be separated cleanly by linear decision boundaries is **linearly separable**.

- **Discriminants.** Classification via $y_k(\mathbf{x})$.

$$\hat{k} = \arg \max_k y_k(\mathbf{x})$$

- **Generative models.** Classification via $p(\mathcal{C}_k)$ and $p(\mathbf{x}|\mathcal{C}_k)$.

$$\hat{k} = \arg \max_k p(\mathcal{C}_k|\mathbf{x}) = \arg \max_k p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$$

- **Discriminative models.** Classification via $p(\mathcal{C}_k|\mathbf{x})$.

$$\hat{k} = \arg \max_k p(\mathcal{C}_k|\mathbf{x})$$

LABELS FOR BINARY CLASSIFICATION

In binary classification, we may use the 0/1 target values for the class labels, i.e. $t = 1$ for a data point of class \mathcal{C}_1 and $t = 0$ for a data point of class \mathcal{C}_2 .

This labeling scheme allows the interpretation that t is the probability of class \mathcal{C}_1 .

In K -ary classification, we often use K **1-hot vectors** for the class labels.

- The 1-hot vector $\mathbf{t}^{(k)}$ has components $t_j^{(k)} = \delta_{kj}$.
- We use $\mathbf{t} = \mathbf{t}^{(k)}$ for a data point of class \mathcal{C}_k .

This labeling scheme allows the interpretation that component t_j is the probability of class \mathcal{C}_j .

Definition. In binary classification, we may make decision based on the posterior probability of class \mathcal{C}_1 . This is achieved by an **activation function** $f(\cdot)$ with range $(0, 1)$ such that

$$p(\mathcal{C}_1|\mathbf{x}) \approx f(a(\mathbf{x}))$$

Note that the activation function act on an **activation**, which is denoted by $a(\mathbf{x})$, of \mathbf{x} .

Definition. A **linear activation** is linear in the input variables

$$a(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

With linear activation, the decision boundary is determined by

$$a(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = f^{-1}\left(\frac{1}{2}\right)$$

Note that this is a hyperplane, and we have a linear model.

FROM INPUT SPACE TO FEATURE SPACE

Consideration in the input space can be extended to a feature space defined by a set of basis functions

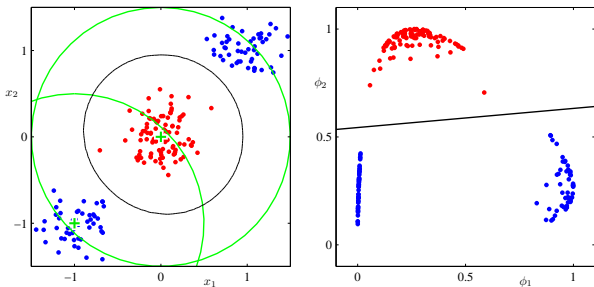
$$\phi = \phi(\mathbf{x})$$

It is often assumed that **activation is linear in the features**. Then we have

$$p(\mathcal{C}_1|\mathbf{x}) \approx f(a(\mathbf{x})) = f(\mathbf{w}^T \phi + w_0)$$

Note that the decision boundary is a hyperplane in a feature space, and non-flat in the input space.

INPUT SPACE AND FEATURE SPACE



The basis functions $\phi_1(\mathbf{x})$ and $\phi_2(\mathbf{x})$ are Gaussian basis functions. A linear decision boundary in the feature space (right) corresponds to a non-linear decision boundary in the input space (left).

In K -ary classification, we may make decision based on the posterior probabilities of K classes. This is achieved by an activation function $\mathbf{y}(\mathbf{x})$ such that the components of \mathbf{y} are non-negative and sum to 1, so

$$(p(\mathcal{C}_1|\mathbf{x}), \dots, p(\mathcal{C}_K|\mathbf{x}))^T \approx \mathbf{f}(\mathbf{a}(\mathbf{x}))$$

Again, the activation function $\mathbf{f}(\cdot)$ acts on class activations

$$\mathbf{a}(\mathbf{x}) = (a_1(\mathbf{x}), \dots, a_K(\mathbf{x}))^T$$

And again, if the class activations are linear functions of the input vector (or feature vector), then we have a (or generalized) linear model.

Discriminants

Definition. A **discriminant** is a function that takes an input vector (or feature vector) and assign it to a class. A discriminant is **linear** if the corresponding decision boundaries are hyperplanes.

SIMPLEST LINEAR DISCRIMINANT

The simplest linear discriminant is a linear function of the input vector

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

Here \mathbf{w} is the **weight vector**, and w_0 is the **bias**. The negative bias is the **threshold**.

Sometimes we write

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

where $\tilde{\mathbf{w}}^T = (w_0, \mathbf{w}^T)$ and $\tilde{\mathbf{x}}^T = (1, \mathbf{x}^T)$.

Definition. In a **binary discriminant**, there is a discriminant function $y(\mathbf{x})$, and an input vector \mathbf{x} is assigned to class \mathcal{C}_1 if $y(\mathbf{x}) \geq 0$ and class \mathcal{C}_2 otherwise. A **binary linear discriminant** has a linear discriminant function.

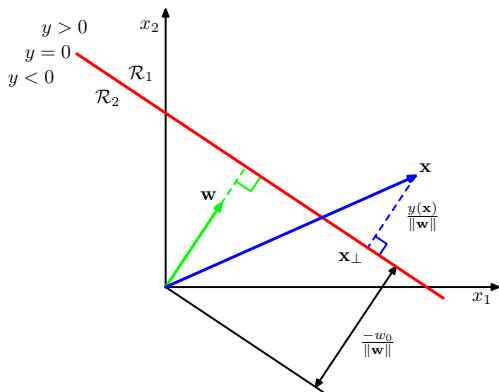
The decision boundary of a binary discriminant is given by

$$y(\mathbf{x}) = 0$$

The decision boundary is a hyperplane for a binary linear discriminant.

PROPERTIES OF A BINARY LINEAR DISCRIMINANT

- w is orthogonal to the decision boundary.
- The distance from x to the decision boundary is $\frac{y(x)}{\|w\|}$.



Definition. In a K -ary discriminant, there are K discriminant functions $y_1(\mathbf{x}), \dots, y_K(\mathbf{x})$, and an input vector \mathbf{x} is assigned to class \mathcal{C}_k if

$$y_k(\mathbf{x}) > y_j(\mathbf{x}) \text{ for all } j \neq k$$

A K -ary linear discriminant has K linear discriminant functions of the input variables

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, \quad k = 1, \dots, K$$

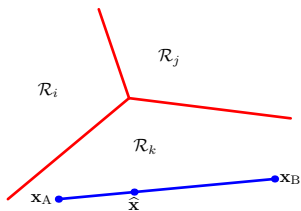
DECISION BOUNDARIES AND DECISION REGIONS

With a K -ary linear discriminant, the decision boundaries are hyperplanes and the decision regions are convex.

The decision boundary between class \mathcal{C}_k and class \mathcal{C}_j is given by $y_k(\mathbf{x}) = y_j(\mathbf{x})$, i.e.

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

which is a hyperplane. Convex decision regions are illustrated below.



LEARNING LINEAR DISCRIMINANTS

The parameters in a linear discriminant (binary or K -ary) can be learned from data.

- Let $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{t}_n)\}_{n=1}^N$ be a data set
- Let $\tilde{\mathbf{X}}$ be the input matrix whose n th row is

$$\tilde{\mathbf{x}}_n^T = (1, x_{n1}, \dots, x_{nD})$$

- Let \mathbf{T} be the target matrix whose n th row is

$$\mathbf{t}_n^T = (t_{n1}, \dots, t_{nK})$$

- Let $\tilde{\mathbf{W}}$ be the weight matrix whose k th column is

$$\tilde{\mathbf{w}}_k = \begin{pmatrix} w_{k0} \\ w_{k1} \\ \vdots \\ w_{kD} \end{pmatrix}$$

MINIMIZING THE SUM OF SQUARED ERRORS

For a K -ary linear discriminant, define the output $\mathbf{y}(\mathbf{x}_n)^T$ of \mathbf{x}_n as the n th row of $\tilde{\mathbf{X}}\tilde{\mathbf{W}}$. The sum of squared errors between the outputs and the targets of \mathcal{D} is

$$\begin{aligned} E(\tilde{\mathbf{W}}) &= \sum_{n=1}^N \|\mathbf{y}(\mathbf{x}_n) - \mathbf{t}_n\|^2 \\ &= \text{tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T \right\} \end{aligned}$$

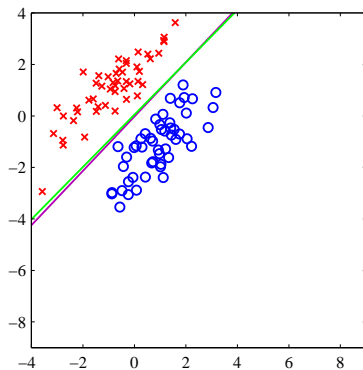
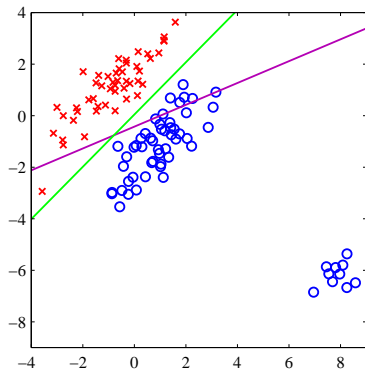
Set the derivative with respect to $\tilde{\mathbf{W}}$ to $\mathbf{0}$

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$$

For a test input vector \mathbf{x} , the output $\mathbf{y}(\mathbf{x})$ is

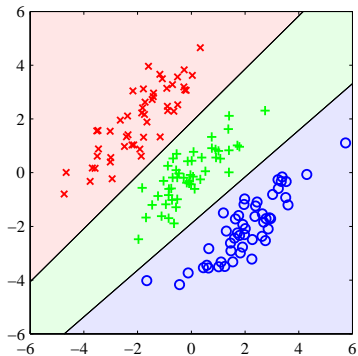
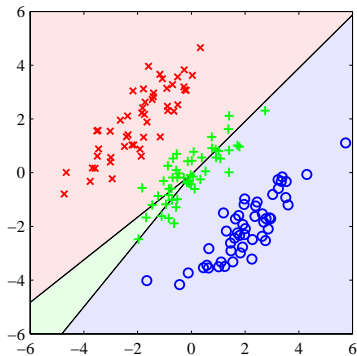
$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} = \mathbf{T}^T \left(\tilde{\mathbf{X}}^\dagger \right)^T \tilde{\mathbf{x}}$$

ISSUE: ROBUSTNESS



Left: With outlier data points
Right: Without outlier data

ISSUE: CORRECTNESS



Left: Least squares learning
Right: Logistic regression learning

Projection is often used in dimension reduction. We want to look for an optimum **direction**.

Consider projection to a line defined by a unit vector \mathbf{w} . Let $m_k \mathbf{w}$ is the mean of the projection points of \mathcal{C}_k , i.e.

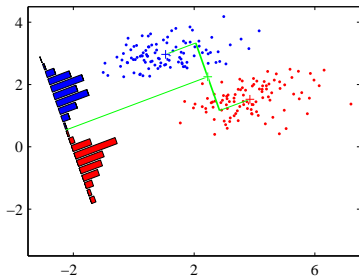
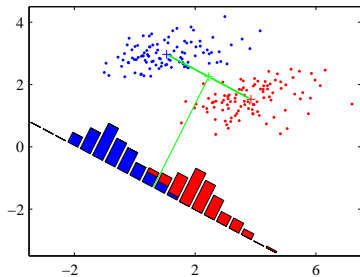
$$m_k \mathbf{w} = \frac{1}{n_k} \sum_{\mathbf{x}_n \in \mathcal{C}_k} (\mathbf{w}^T \mathbf{x}_n) \mathbf{w} = \mathbf{w}^T \left(\frac{1}{n_k} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n \right) \mathbf{w} = (\mathbf{w}^T \mathbf{m}_k) \mathbf{w}$$

where $\mathbf{m}_k = \frac{1}{n_k} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n$ is the mean of data points of class \mathcal{C}_k . The separation between classes \mathcal{C}_1 and \mathcal{C}_2 after the projection can be measured by $(m_2 - m_1)^2$. Since $(m_2 - m_1)^2 = (\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1))^2$, it is maximized when

$$\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$$

As is shown, there may be considerable overlap between the classes when

$$w \propto (m_2 - m_1)$$



Idea. Minimize the class overlap by large separation between class means and small variance within each class.

Let \mathbf{w} be a unit vector defining a line and $y_n = \mathbf{w}^T \mathbf{x}_n$, so $y_n \mathbf{w}$ is the projection of \mathbf{x}_n on the line. We define the **within-class variance** of class \mathcal{C}_k as

$$s_k^2 = \sum_{\mathbf{x}_n \in \mathcal{C}_k} (y_n - m_k)^2$$

and the **total within-class variance** as $s_1^2 + s_2^2$.

Criterion. Maximize

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

COVARIANCE MATRICES BASED ON DATA

The numerator of $J(\mathbf{w})$ can be written as

$$(m_2 - m_1)^2 = \mathbf{w}^T (m_2 - m_1)(m_2 - m_1)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_B \mathbf{w}$$

where \mathbf{S}_B is the **between-class covariance matrix** defined by

$$\mathbf{S}_B = (m_2 - m_1)(m_2 - m_1)^T$$

The denominator can be written as

$$s_1^2 + s_2^2 = \sum_{\mathbf{x}_n \in \mathcal{C}_1} (y_n - m_1)^2 + \sum_{\mathbf{x}_n \in \mathcal{C}_2} (y_n - m_2)^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w}$$

where \mathbf{S}_W is the **total within-class covariance matrix** defined by

$$\mathbf{S}_W = \sum_{\mathbf{x}_n \in \mathcal{C}_1} (\mathbf{x}_n - m_1)(\mathbf{x}_n - m_1)^T + \sum_{\mathbf{x}_n \in \mathcal{C}_2} (\mathbf{x}_n - m_2)(\mathbf{x}_n - m_2)^T$$

Thus, the Fisher's criterion can be written as

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

Setting the gradient of $J(\mathbf{w})$ to $\mathbf{0}$, we get

$$\begin{aligned} (\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} &= (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w} \\ &\propto \mathbf{S}_B \mathbf{w} \\ &\propto (\mathbf{m}_2 - \mathbf{m}_1) \end{aligned}$$

Hence

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

For K classes, we define the **between-class covariance matrix** as

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

the **total within-class covariance matrix** as

$$\mathbf{S}_W = \sum_{k=1}^K \sum_{\mathbf{x}_n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$$

and the **total covariance matrix** as

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m})(\mathbf{x}_n - \mathbf{m})^T$$

Note that $\mathbf{S}_B + \mathbf{S}_W = \mathbf{S}_T$.

MULTIPLE LINEAR FEATURES

For K classes, more than one directions are needed.

Through w , a **linear feature** $y(x) = w^T x$ is extracted from x . Multiple linear features can be extracted through multiple w 's

$$y_i(x) = w_i^T x, \quad i = 1, \dots, D'$$

These features form a **feature vector**

$$\mathbf{y}(x) = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{D'} \end{bmatrix} = \begin{bmatrix} w_1^T x \\ w_2^T x \\ \vdots \\ w_{D'}^T x \end{bmatrix} = \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_{D'}^T \end{bmatrix} x = \mathbf{W}^T x$$

The weight matrix \mathbf{W} has the weight vectors $w_1, \dots, w_{D'}$ as columns.

COVARIANCE MATRICES

Define $\mathbf{y}_n = \mathbf{W}^T \mathbf{x}_n$ and

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{y}_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k$$

The total within-class covariance matrix in the feature space is

$$\begin{aligned} \mathbf{s}_W &= \sum_{k=1}^K \sum_{\mathbf{x}_n \in \mathcal{C}_k} (\mathbf{y}_n - \boldsymbol{\mu}_k)(\mathbf{y}_n - \boldsymbol{\mu}_k)^T \\ &= \mathbf{W}^T \left(\sum_{k=1}^K \sum_{\mathbf{x}_n \in \mathcal{C}_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T \right) \mathbf{W} = \mathbf{W}^T \mathbf{S}_W \mathbf{W} \end{aligned}$$

The between-class covariance matrix is

$$\mathbf{s}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu})(\boldsymbol{\mu}_k - \boldsymbol{\mu})^T = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$$

FISHER'S CRITERION FOR MULTIPLE FEATURES

$$\begin{aligned} J(\mathbf{W}) &= \text{tr}\{\mathbf{s}_W^{-1} \mathbf{s}_B\} \\ &= \text{tr}\{(\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{S}_B \mathbf{W})\} \end{aligned}$$

- $J(\mathbf{W})$ is large when the between-class covariance is large and when the within-class covariance is small.
- \mathbf{W} is determined by the eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$ with the D' largest eigenvalues.

Perceptron

CLASS LABELS IN PERCEPTRON

In perceptron, we use class label $t \in \{-1, 1\}$, where $t = 1$ represents \mathcal{C}_1 and $t = -1$ represents \mathcal{C}_2 , i.e.

$$t_n = \begin{cases} 1, & \text{if } \mathbf{x}_n \in \mathcal{C}_1 \\ -1, & \text{if } \mathbf{x}_n \in \mathcal{C}_2 \end{cases}$$

and the decision function of

$$y(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \boldsymbol{\phi})$$

where $\boldsymbol{\phi} = \boldsymbol{\phi}(\mathbf{x})$ is the feature vector based on fixed basis functions.

The parameters \mathbf{w} are learned by a set $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$.

A perceptron makes an error for a data point (\mathbf{x}, t) if

$$\mathbf{w}^T \phi t < 0$$

The **error measure** of data point (\mathbf{x}_n, t_n) can be defined by

$$E_n(\mathbf{w}) = \begin{cases} -\mathbf{w}^T \phi_n t_n, & \mathbf{w}^T \phi_n t_n < 0 \\ 0, & \text{otherwise} \end{cases}$$

For a misclassified example (\mathbf{x}_n, t_n) , update \mathbf{w} by

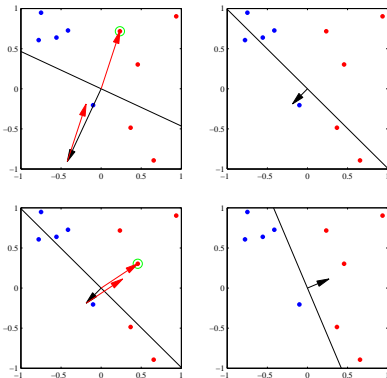
$$\begin{aligned}\mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) \\ &= \mathbf{w}^{(\tau)} + \eta \phi_n t_n\end{aligned}$$

where $\eta > 0$ is a learning rate.

- For misclassified \mathbf{x}_n with $t_n = -1$, subtract $\eta \phi_n$ from \mathbf{w} .
- For misclassified \mathbf{x}_n with $t_n = 1$, add $\eta \phi_n$ from \mathbf{w} .

PERCEPTRON CONVERGENCE THEOREM

For linearly separable data, the perceptron algorithm learns a perfect decision boundary in finite steps.



Probabilistic Generative Model

Definition. A **generative model** consists of the **class priors** $p(\mathcal{C}_k)$ and the **class-conditional distributions** $p(\mathbf{x}|\mathcal{C}_k)$.

- **Data generation.** A generative model can be used to generate artificial data.
- **Class posteriors.** By Bayes' rule, the posterior of class \mathcal{C}_k is

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}, \mathcal{C}_k)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \\ &\propto p(\mathbf{x}|\mathcal{C}_k) p(\mathcal{C}_k) \end{aligned}$$

$p(\mathcal{C}_1|\mathbf{x}), \dots, p(\mathcal{C}_K|\mathbf{x})$ are used in making decisions.

In binary classification of classes \mathcal{C}_1 and \mathcal{C}_2 , the posterior of \mathcal{C}_1 is

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}} \\ &= \frac{1}{1 + \exp(-a)} \\ &= \sigma(a) \end{aligned}$$

where $\exp(-a) = \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}$ or $a = \log \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$.

Definition. The function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

is the **logistic sigmoid function**.

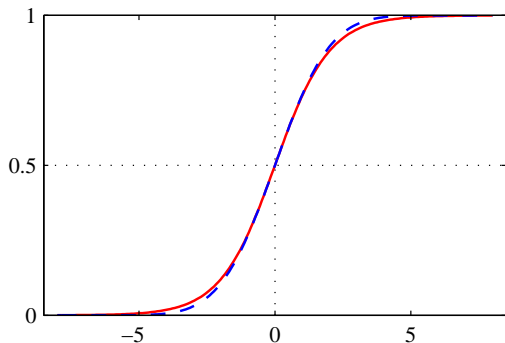
- S-shaped
- squashing: mapping $(-\infty, \infty)$ to $(0, 1)$
- symmetry property

$$\sigma(-a) = 1 - \sigma(a)$$

- a differentiable approximation of the step function with

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

SIGMOID AND INVERSE PROBIT FUNCTIONS



Plot of the logistic sigmoid function $\sigma(a)$ (as shown in red) and a scaled inverse probit function $\Phi(\lambda a)$ for $\lambda = \sqrt{\frac{\pi}{8}}$ (blue)

Definition. The **logit function** is the inverse of the logistic sigmoid function.

$$a(\sigma) = \log \left(\frac{\sigma}{1 - \sigma} \right)$$

In binary classification of classes \mathcal{C}_1 and \mathcal{C}_2 , we often have $\sigma = p(\mathcal{C}_1|\mathbf{x})$. Then the logit

$$a = \log \left(\frac{\sigma}{1 - \sigma} \right) = \log \left(\frac{p(\mathcal{C}_1|\mathbf{x})}{1 - p(\mathcal{C}_1|\mathbf{x})} \right) = \log \left(\frac{p(\mathcal{C}_1|\mathbf{x})}{p(\mathcal{C}_2|\mathbf{x})} \right)$$

is the **log odds**.

In K -ary classification of classes $\mathcal{C}_1, \dots, \mathcal{C}_K$, the posterior of \mathcal{C}_k is

$$\begin{aligned} p(\mathcal{C}_k|\mathbf{x}) &= \frac{p(\mathbf{x}, \mathcal{C}_k)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

where $\exp(a_k) = p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ or $a_k = \log(p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k))$. We call a_k the **activation** of class \mathcal{C}_k .

Definition. A common transformation from class actions to class posterior probabilities is the **normalized exponential function** defined by

$$(a_1, \dots, a_K) \mapsto \frac{(\exp(a_1), \dots, \exp(a_K))}{\sum_{j=1}^K \exp(a_j)}$$

The normalized exponential is also known as the **softmax**: for the class \mathcal{C}_k with the largest a_k , we have

$$\frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)} \approx 1, \quad \frac{\exp(a_{j \neq k})}{\sum_{j=1}^K \exp(a_j)} \approx 0$$

if $a_k \gg a_j$ for $j \neq k$.

In a generative model

- The class priors are simply the relative frequencies.
- The key elements are the class-conditional distributions.
- We now introduce a few instances.

Definition. A **Gaussian class-conditional distribution** is

$$\begin{aligned} p(\mathbf{x}|\mathcal{C}_k) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}-\boldsymbol{\mu}_k)} \end{aligned}$$

where $\boldsymbol{\mu}_k$ is class mean and $\boldsymbol{\Sigma}_k$ is class covariance matrix.

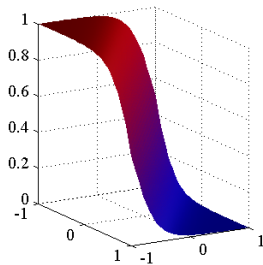
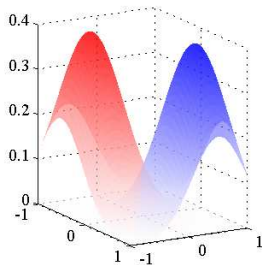
A **Gaussian class-conditional model** assumes a Gaussian class-conditional distribution for each class.

Consider a generative model with Gaussian class-conditional distributions.

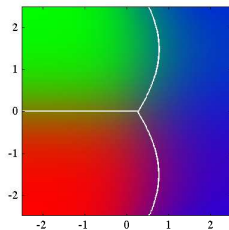
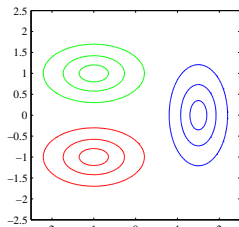
- If the Gaussian class-conditional distributions do not share a common covariance matrix, then the decision boundaries are quadratic functions of the input vector.
- If the Gaussian class-conditional distributions share a common covariance matrix, i.e.

$$\begin{aligned} p(\mathbf{x}|\mathcal{C}_k) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}) \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_k)} \end{aligned}$$

then the decision boundaries are hyperplanes.



Class-conditionals $p(\mathbf{x}|\mathcal{C}_k)$ (left) and posterior $p(\mathcal{C}_1|\mathbf{x})$ (right).



Class-conditionals (left) and decision boundaries (right).

Substituting the Gaussian class-conditional distributions and eliminating common factors (the quadratic terms), we get

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0) = \sigma(a)$$

where

$$\begin{aligned} \mathbf{w} &= \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ w_0 &= -\frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \log \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)} \end{aligned}$$

The decision boundary is the hyperplane $a = \mathbf{w}^T \mathbf{x} + w_0 = 0$.

Substituting the Gaussian class-conditional distributions and eliminating common factors (the quadratic terms), we get

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{\exp(\mathbf{w}_k^T \mathbf{x} + w_{k0})}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x} + w_{j0})} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k, \quad w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log p(\mathcal{C}_k)$$

The boundary between \mathcal{R}_j and \mathcal{R}_k is decided by $a_j = a_k$, i.e.

$$\mathbf{w}_j^T \mathbf{x} + w_{j0} = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

A **naïve Bayes model** is a generative model assuming independent features for the class-conditional distributions.

In particular, for **binary features** $\mathbf{x} \in \{0, 1\}^D$, the class-conditional distributions are

$$p(\mathbf{x}|\mathcal{C}_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}, \quad k = 1, \dots, K$$

A naïve Bayes model with binary features is a linear model for classification. The activations are linear in \mathbf{x} , as

$$\begin{aligned} a_k &= \log(p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)) \\ &= \sum_{i=1}^D \{x_i \log \mu_{ki} + (1 - x_i) \log(1 - \mu_{ki})\} + \log p(\mathcal{C}_k) \end{aligned}$$

Discriminative Model

- In a generative model, we compute the posterior probabilities based on the prior probabilities and conditional distributions.
- In a discriminative model, we model the class posterior probabilities directly.
- In particular, in a **generalized linear discriminative model**, we assume that the class activations are linear functions of the feature vector

$$a_k = \mathbf{w}_k^T \boldsymbol{\phi} + w_{k0}$$

- Again, the activations are transformed to the posterior probabilities by a non-linear function, e.g. logistic sigmoid or normalized exponential.

Definition. A **logistic regression** models the posterior probability of class \mathcal{C}_1 as the **logistic sigmoid function** acting on a linear function of the feature vector

$$p(\mathcal{C}_1|\mathbf{x}) \approx y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \phi)$$

The posterior probability of class \mathcal{C}_2 is

$$p(\mathcal{C}_2|\mathbf{x}) = 1 - p(\mathcal{C}_1|\mathbf{x}) = 1 - \sigma(\mathbf{w}^T \phi)$$

The parameters \mathbf{w} in a logistic regression can be learned by a data set $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$.

The data likelihood of (\mathbf{x}, t) is $p(\mathcal{C}_1|\mathbf{x})$ if $t = 1$ and $p(\mathcal{C}_2|\mathbf{x})$ if $t = 0$. Since $p(\mathcal{C}_1|\mathbf{x}) = y(\mathbf{x}, \mathbf{w})$ and $p(\mathcal{C}_2|\mathbf{x}) = 1 - y(\mathbf{x}, \mathbf{w})$, the data likelihood of (\mathbf{x}, t) can be written as

$$p(t|\mathbf{x}) = y^t(1 - y)^{1-t}$$

where $y = y(\mathbf{x}, \mathbf{w})$. The likelihood of \mathcal{D} is

$$p(\mathcal{D}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

where $y_n = y(\mathbf{x}_n, \mathbf{w})$.

The negative log likelihood function (error function) is

$$\begin{aligned} E(\mathbf{w}) &= -\log p(\mathcal{D}|\mathbf{w}) \\ &= -\log \left(\prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n} \right) \\ &= -\sum_{n=1}^N \{t_n \log y_n + (1 - t_n) \log(1 - y_n)\} \\ &= \sum_{n=1}^N E_n(\mathbf{w}) \end{aligned}$$

Note that $E_n(\mathbf{w}) = -\{t_n \log y_n + (1 - t_n) \log(1 - y_n)\}$ is the **cross entropy** between $(t_n, 1 - t_n)$ and $(y_n, 1 - y_n)$.

GRADIENT OF ERROR FUNCTION

Recall that $y(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \boldsymbol{\phi})$. The gradient of $E(\mathbf{w})$ is

$$\begin{aligned}\nabla E(\mathbf{w}) &= - \sum_{n=1}^N \nabla \{t_n \log y_n + (1 - t_n) \log(1 - y_n)\} \\ &= - \sum_{n=1}^N \frac{t_n}{y_n} \nabla y_n + \left(\frac{1 - t_n}{1 - y_n} \right) \nabla(1 - y_n) \\ &= - \sum_{n=1}^N \frac{t_n}{y_n} y_n (1 - y_n) \boldsymbol{\phi}_n - \left(\frac{1 - t_n}{1 - y_n} \right) y_n (1 - y_n) \boldsymbol{\phi}_n \\ &= \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n \\ &= \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})\end{aligned}$$

where $\boldsymbol{\Phi}$ is the design matrix.

In **gradient decent**, the parameters are updated by

$$\begin{aligned}\mathbf{w}' &= \mathbf{w} - \eta \nabla E(\mathbf{w}) \\ &= \mathbf{w} - \eta \Phi^T(\mathbf{y} - \mathbf{t})\end{aligned}$$

In **stochastic gradient decent**, \mathbf{w} are updated by

$$\begin{aligned}\mathbf{w}' &= \mathbf{w} - \eta \nabla E_n(\mathbf{w}) \\ &= \mathbf{w} - \eta(y_n - t_n)\phi_n\end{aligned}$$

where (t_n, \mathbf{x}_n) is an example in \mathcal{D} .

NEWTON-RAPHSON METHOD

In **Newton-Raphson method**, w are updated by

$$w' = w - H^{-1} \nabla E(w)$$

where H is the Hessian matrix. For the cross entropy error function

$$H = \nabla \Phi^T (\mathbf{y} - \mathbf{t}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

where $\mathbf{R} = \text{diag}\{y_1(1 - y_1), \dots, y_n(1 - y_n)\}$. Thus

$$\begin{aligned} w' &= w - H^{-1} \nabla E(w) = w - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi w - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} z \end{aligned}$$

where $z = \Phi w - \mathbf{R}^{-1}(\mathbf{y} - \mathbf{t})$.

WEIGHTED LEAST-SQUARES PROBLEMS

A solution of the **least-squares problem** $Ax = b$ satisfies the normal equation

$$A^T Ax = A^T b$$

A solution of the **weighted least-squares problem** $C Ax = C b$, where C is diagonal with positive diagonal elements, satisfies the normal equation

$$A^T C^T C A x = A^T C^T C b$$

That is

$$x = \left(A^T C^T C A \right)^{-1} A^T C^T C b$$

Recall that the weight is updated by

$$\mathbf{w}' = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}$$

This is the solution of a weighted least-squares problem

$$\mathbf{C} \Phi \mathbf{w} = \mathbf{C} \mathbf{z}$$

where

$$\mathbf{R} = \mathbf{C}^T \mathbf{C}$$

Since \mathbf{R} depends on \mathbf{w} , the weights are updated in each iteration.

Definition. A **multiclass logistic regression** models the class posterior probabilities as the normalized exponential function acting on linear functions of the feature vector. Specifically, the posterior probability of class \mathcal{C}_k

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{\exp(a_k)}{\sum_{j=1}^K \exp(a_j)}$$

where the class **activations** are given by

$$a_k = \mathbf{w}_k^T \phi$$

The parameters $\mathbf{w}_1, \dots, \mathbf{w}_K$ in a multiclass logistic regression can be learned from data.

The likelihood of an example (\mathbf{x}, \mathbf{t}) is

$$p(\mathbf{t}|\mathbf{x}) = \prod_{k=1}^K y_k^{t_k}$$

where

$$y_k = y_k(\mathbf{x}) = \frac{\exp(a_k(\mathbf{x}))}{\sum_{j=1}^K \exp(a_j(\mathbf{x}))}$$

The likelihood of a data set $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{t}_n)\}_{n=1}^N$ is

$$p(\mathcal{D}|\mathbf{w}) = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

CROSS-ENTROPY ERROR FUNCTION

Maximizing the data likelihood is equivalent to minimizing the negative log data likelihood, so we define

$$\begin{aligned} E(\mathbf{w}) &= -\log \prod_{n=1}^N p(\mathbf{t}_n | \mathbf{x}_n) \\ &= -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk} \\ &= \sum_{n=1}^N E_n(\mathbf{w}) \end{aligned}$$

Note that

$$E_n = -\sum_{k=1}^K t_{nk} \log y_{nk}$$

is the cross entropy between \mathbf{t}_n and \mathbf{y}_n .

From $a_j = \mathbf{w}_j^T \boldsymbol{\phi}$ and $y_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$, we have

$$\nabla_{\mathbf{w}_j} a_l = \delta_{jl} \boldsymbol{\phi} \quad \text{and} \quad \frac{\partial y_k}{\partial a_j} = y_k (\delta_{kj} - y_j)$$

It follows that

$$\begin{aligned} \nabla_{\mathbf{w}_j} y_k &= \sum_l \frac{\partial y_k}{\partial a_l} (\nabla_{\mathbf{w}_j} a_l) \\ &= \sum_l y_k (\delta_{kl} - y_l) \delta_{jl} \boldsymbol{\phi} \\ &= y_k (\delta_{kj} - y_j) \boldsymbol{\phi} \end{aligned}$$

Consider the first-order derivatives of the error $E = -\sum_{k=1}^K t_k \log y_k$ of a data point (\mathbf{t}, \mathbf{x}) .

$$\begin{aligned}
 \nabla_{\mathbf{w}_j} \left(-\sum_{k=1}^K t_k \log y_k \right) &= -\sum_{k=1}^K t_k \left(\frac{1}{y_k} \nabla_{\mathbf{w}_j} y_k \right) \\
 &= -\sum_{k=1}^K t_k ((\delta_{kj} - y_j) \phi) \\
 &= -t_j \phi + \sum_{k=1}^K t_k y_j \phi \\
 &= -t_j \phi + y_j \phi \\
 &= (y_j - t_j) \phi
 \end{aligned}$$

Thus, the gradient of $E(\mathbf{w})$ of data set \mathcal{D} is

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}) = \nabla_{\mathbf{w}_j} \left(\sum_{n=1}^N E_n(\mathbf{w}) \right) = \sum_{n=1}^N (y_{nj} - t_{nj}) \phi_n = \mathbf{\Phi}^T (\mathbf{y}_j - \mathbf{t}_j)$$

Consider the second-order derivatives of the error $E = -\sum_{k=1}^K t_k \log y_k$ of a data point (\mathbf{t}, \mathbf{x}) .

$$\begin{aligned}\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} \left(-\sum_{k=1}^K t_k \log y_k \right) &= \nabla_{\mathbf{w}_k} (y_j - t_j) \phi \\ &= y_k (\delta_{kj} - y_j) \phi \phi^T\end{aligned}$$

Thus, the Hessian of $E(\mathbf{w})$ of data set \mathcal{D} is

$$\begin{aligned}\nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}) &= \nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} \left(\sum_{n=1}^N E_n(\mathbf{w}) \right) = \sum_{n=1}^N y_{nk} (\delta_{kj} - y_{nj}) \phi_n \phi_n^T \\ &= \mathbf{\Phi}^T \mathbf{R}' \mathbf{\Phi}\end{aligned}$$

The Newton-Raphson method (via iterative re-weighted least squares) can be used to learn parameter iteratively.

Other than the logistic sigmoid function, there are many choices for the activation function in a binary classification.

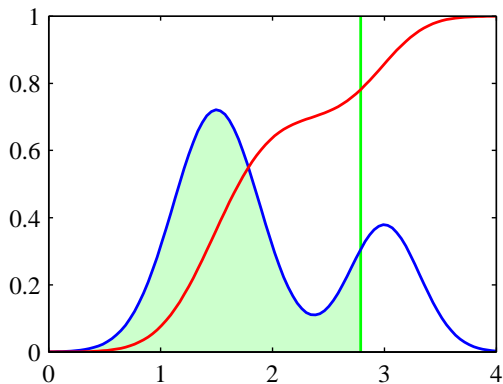
In fact, any CDF

$$f(a) = \int_{-\infty}^a p(t) dt$$

where $p(\cdot)$ is a PDF, is a valid activation function for binary classification since

$$0 \leq f(a) \leq 1$$

CDF AND PDF



Definition. The CDF of a standard Gaussian

$$\Phi(a) = \int_{-\infty}^a \mathcal{N}(t|0, 1)dt$$

is the **inverse probit function**. A generalized linear model with the inverse probit function as the activation function is a **probit regression**.

The data labels may be subject to errors.

For binary classification, the conditional likelihood of (t, \mathbf{x}) is

$$\begin{aligned} p(t|\mathbf{x}) &= \sum_{e=0}^1 p(t, e|\mathbf{x}) \\ &= \sum_{e=0}^1 p(e)p(t|e, \mathbf{x}) \\ &= \begin{cases} (1 - \epsilon)\sigma(\mathbf{x}) + \epsilon(1 - \sigma(\mathbf{x})), & t = 1 \\ \epsilon\sigma(\mathbf{x}) + (1 - \epsilon)(1 - \sigma(\mathbf{x})), & t = 0 \end{cases} \\ &= [(1 - \epsilon)\sigma(\mathbf{x}) + \epsilon(1 - \sigma(\mathbf{x}))]^t [\epsilon\sigma(\mathbf{x}) + (1 - \epsilon)(1 - \sigma(\mathbf{x}))]^{1-t} \end{aligned}$$

where $\epsilon = p(e)$ is the probability of error in labeling.

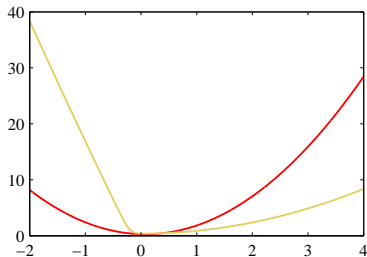
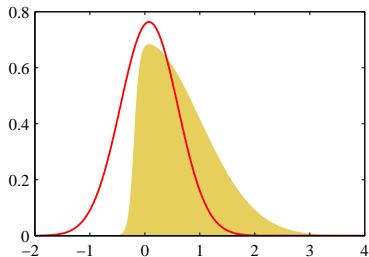
Laplace Approximation

APPROXIMATE PDF BY GAUSSIAN

Let $p(\mathbf{x})$ be a PDF defined by unnormalized $f(\mathbf{x}) \geq 0$

$$p(\mathbf{x}) = \frac{1}{Z} f(\mathbf{x})$$

In Laplace approximation, $p(\mathbf{x})$ is approximated by a Gaussian PDF centered on a mode x_0 of $f(\mathbf{x})$ (or $p(\mathbf{x})$). The precision of the Gaussian is the negative Hessian of $\log f(\mathbf{x})$ at x_0 .



UNI-VARIATE LAPLACE APPROXIMATION

- Given $f(z)$, find a **mode** (local maximum) z_0 of $f(z)$ with

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

- Find the negative second derivative of $\log f(z)$ at z_0

$$A = - \left. \frac{d^2}{dz^2} \log f(z) \right|_{z=z_0}$$

- The second-order approximation of $\log f(z)$ around z_0 is

$$\log f(z) \approx \log f(z_0) - \frac{1}{2} A (z - z_0)^2$$

- The Laplace approximation of $p(z)$ is

$$\begin{aligned} q(z) &\propto \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\} \\ &= \left(\frac{A}{2\pi} \right)^{1/2} \exp \left\{ -\frac{A}{2} (z - z_0)^2 \right\} \end{aligned}$$

MULTI-VARIATE LAPLACE APPROXIMATION

- Given $f(\mathbf{z})$, find a **mode** (local maximum) \mathbf{z}_0 of $f(\mathbf{z})$ with

$$\nabla f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0} = \mathbf{0}$$

- Find the negative Hessian of $\log f(\mathbf{z})$ at \mathbf{z}_0

$$\mathbf{A} = -\nabla\nabla \log f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0}$$

- The second-order approximation of $\log f(\mathbf{z})$ around \mathbf{z}_0 is

$$\log f(\mathbf{z}) \approx \log f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)$$

- The Laplace approximation of $p(\mathbf{z})$ is

$$\begin{aligned} q(\mathbf{z}) &\propto \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} \\ &= \mathcal{N}(\mathbf{z} | \mathbf{z}_0, \mathbf{A}^{-1}) \end{aligned}$$

Bayesian Logistic Regression

Consider a Gaussian prior for the parameters

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

For the posterior of the parameters, we have

$$p(\mathbf{w} | \mathcal{D}) \propto p(\mathbf{w}) p(\mathcal{D} | \mathbf{w})$$

Taking the logarithm, we have

$$\begin{aligned} \log p(\mathbf{w} | \mathcal{D}) = & -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\ & + \sum_{n=1}^N t_n \log y_n + (1 - t_n) \log(1 - y_n) + \text{const} \end{aligned}$$

LAPLACE APPROXIMATION

To proceed, we apply the Laplace approximation to $p(\mathbf{w}|\mathcal{D})$.

- The mode of $p(\mathbf{w}|\mathcal{D})$ is \mathbf{w}_{MAP}
- The negative Hessian of $\log p(\mathbf{w}|\mathcal{D})$ at \mathbf{w}_{MAP} is

$$-\nabla\nabla \log p(\mathbf{w}|\mathcal{D}) = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n)\phi_n\phi_n^T$$

Hence, the Laplace approximation to $p(\mathbf{w}|\mathcal{D})$ is

$$p(\mathbf{w}|\mathcal{D}) \approx q(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

where

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \sum_{n=1}^N y_n(1 - y_n)\phi_n\phi_n^T$$

The prediction is the integration of the prediction over the posterior distribution of the parameters.

With the Laplace approximation, we have

$$\begin{aligned} p(C_1|\mathbf{x}, \mathcal{D}) &= \int p(C_1|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w} \\ &\approx \int \sigma(\mathbf{w}^T \boldsymbol{\phi})q(\mathbf{w})d\mathbf{w} \\ &= \int \sigma(a)p(a)da \end{aligned}$$

where $a = \mathbf{w}^T \boldsymbol{\phi}$.

Note that $a = \mathbf{w}^T \boldsymbol{\phi}$ is Gaussian as \mathbf{w} is Gaussian. In particular, the mean and variance of $a = \mathbf{w}^T \boldsymbol{\phi}$ is

$$\mu_a = \int a(\mathbf{w}) q(\mathbf{w}) d\mathbf{w} = \int \mathbf{w}^T \boldsymbol{\phi} q(\mathbf{w}) d\mathbf{w} = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}$$

$$\sigma_a^2 = \int (\mathbf{w}^T \boldsymbol{\phi} - \mu_a)^2 q(\mathbf{w}) d\mathbf{w} = \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}$$

Thus

$$\int \sigma(a) p(a) da = \int \sigma(a) \mathcal{N}(a | \mu_a, \sigma_a^2) da$$

FURTHER APPROXIMATION

A logistic sigmoid function can be approximated by an inverse probit function

$$\sigma(a) \approx \Phi(\lambda a), \quad \lambda^2 = \frac{\pi}{8}$$

It can be shown that

$$\int \Phi(\lambda a) \mathcal{N}(a|\mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$

Thus we have

$$\int \sigma(a) \mathcal{N}(a|\mu, \sigma^2) da \approx \sigma\left(\kappa(\sigma^2)\mu\right), \quad \kappa(\sigma^2) = (1 + \lambda^2\sigma^2)^{-1/2}$$

So

$$p(\mathcal{C}_1|\mathbf{x}, \mathcal{D}) \approx \sigma\left(\kappa(\sigma_a^2)\mu_a\right)$$

Note that the decision boundary is given by $\mu_a = \mathbf{w}_{\text{MAP}}^T \phi = 0$.