

5. Logistic Regression for Classification

Dongwoo Kim

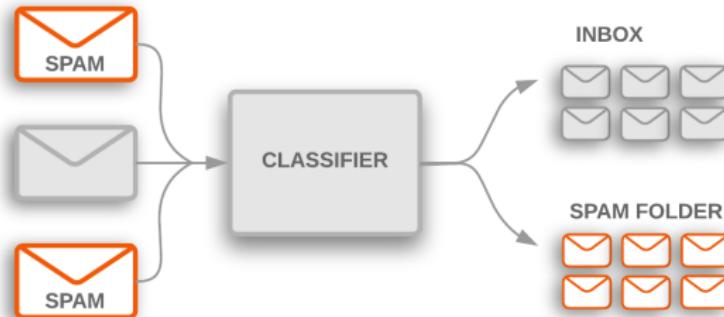
dongwookim.ac.kr

CSED515 - 2023 Spring

Classification

Goal: to assign an input feature vector x to one of K discrete class \mathcal{C}_k where $k = 1, \dots, K$.

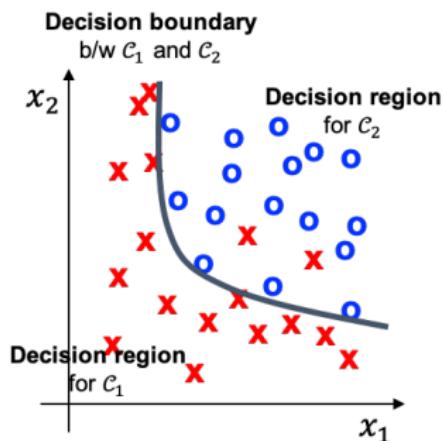
- e.g., handwriting recognition, speech recognition, biological classification, drug discovery, ..., **spam filtering**



<https://medium.com/@naveen.kumar.k/naive-bayes-spam-detection-7d087cc96d9d>

Decision Boundary

- Suppose that there are two classes $\mathcal{C}_1, \mathcal{C}_2$
- A classifier assigns a new data point x to \mathcal{C}_1 only if it is in the decision region \mathcal{R}_1 for \mathcal{C}_1 , i.e., classification function $h(x) = k$ if $x \in \mathcal{R}_k$.
- Decision boundaries (or surface) are boundaries between decision regions



developing classifier = drawing decision boundaries

Linear Regression for Classification?

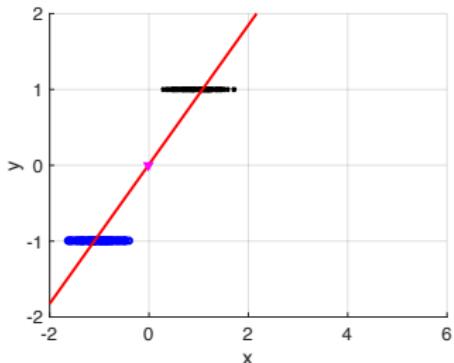
- ▶ Consider dataset $\mathcal{D} = (x_n, y_n)_{n \in [N]}$ for classification where $x_n \in \mathbb{R}$ and $y_n = \{-1, +1\}$.
- ▶ Assuming linear model $y = w_1 x + w_2$, then a classifier h can be defined as:

$$h(x) = \text{sign}(w_1 x + w_2) = \begin{cases} +1 & \text{if } w_1 x + w_2 \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

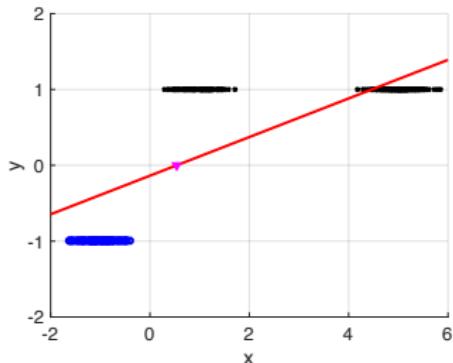
- ▶ Learning the linear regression aims at finding $w = (w_1, w_2)$ such that the **square error** is small

$$\arg \min_w \frac{1}{2} \sum_{n=1}^N (y_n - (w_1 x_n + w_2))^2$$

Linear Regression for Classification?



perfect classification



decision boundary shifted

- ▶ Note that x_i having $w_1x_i + w_2$ farther from zero is easier sample
- ▶ Employing square error penalizes samples that are **very easy to classify**
- ▶ This problem is from the large **mismatch** between linear model $y = w_2x + w_1$ and data $(x_n, y_n)_{n=1,\dots,N}$

Outline

- 1 A systemic formulation: logistic regression
 - The optimality of Bayes decision rule
 - Logistic regression and probabilistic model
- 2 A systemic solver: gradient-based optimization
 - Gradient decent/ascent
 - Newton method

Discriminant Function

An alternative description of classifier:

assign x to \mathcal{C}_k if $f_k(x) > f_j(x)$ for all $j \neq k$

- ▶ $\{f_k\}_{k=1,\dots,K}$: discriminant function/score
 - ▶ e.g., $f_k(x) = p(\mathcal{C}_k | x)$
- ▶ i.e., classifier $h(x) = \arg \max_k f_k(x)$
- ▶ If $f_k(\cdot)$'s are a set of discriminant functions, then for any monotonically increasing function $g : \mathbb{R} \rightarrow \mathbb{R}$ $g_k(\cdot) = g(f_k(\cdot))$'s are also a set of discriminant functions
 - ▶ e.g.1, $a f_k(x) + b$ for any $a > 0$ and b
 - ▶ e.g.2, $\log f_k(\cdot)$
 - ▶ ...

Discriminant Function: Bayes Decision Rule

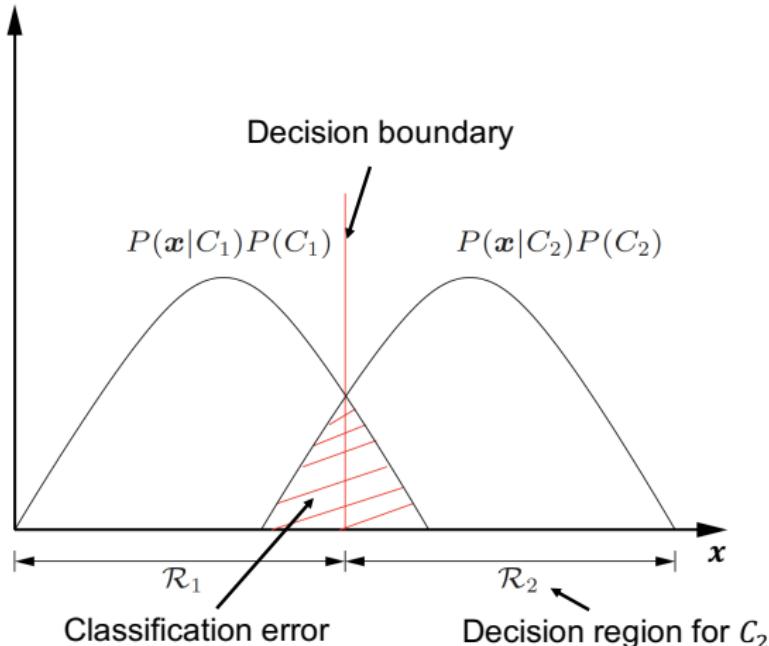
Bayes decision rule:

$$f_k(x) = p(\mathcal{C}_k | x) \text{ or equivalently, } f_k(x) = p(x | \mathcal{C}_k)p(\mathcal{C}_k)$$

- Classify $x \in \mathcal{C}_k$ if $p(\mathcal{C}_k | x) > p(\mathcal{C}_j | x)$ for all $j \neq k$.
- We can use discriminant function $f_k(x) = p(x | \mathcal{C}_k)p(\mathcal{C}_k)$ without **normalization factor** thanks to Bayes rule:

$$p(\mathcal{C}_k | x) = \frac{\overbrace{p(x | \mathcal{C}_k)}^{\text{Class-conditional density}} \overbrace{p(\mathcal{C}_k)}^{\text{Prior}}}{\underbrace{\sum_{j=1}^K p(x | \mathcal{C}_j)p(\mathcal{C}_j)}_{\text{Normalization factor}}}$$

Graphical Illustration of Bayes Decision Rule



Optimality of Bayes Decision Rule

The accuracy or $(1 - p(\text{error}))$ is given by:

$$\begin{aligned} p(\text{correct}) &= \sum_{k=1}^K p(x \in \mathcal{R}_k, \mathcal{C}_k) \\ &= \sum_{k=1}^K p(x \in \mathcal{R}_k \mid \mathcal{C}_k) p(\mathcal{C}_k) \\ &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(x \mid \mathcal{C}_k) p(\mathcal{C}_k) dx . \end{aligned}$$

By dividing the domain of x into infinitesimal regions ε 's, in order to maximize the accuracy (or to minimize the error rate), a region ε should be included in \mathcal{R}_k such that $p(x \mid \mathcal{C}_k)p(\mathcal{C}_k) > p(x \mid \mathcal{C}_j)p(\mathcal{C}_j)$ for all $j \neq k$. (A formal argument can be found at [link](#))

Outline

- 1 A systemic formulation: logistic regression
 - The optimality of Bayes decision rule
 - Logistic regression and probabilistic model
- 2 A systemic solver: gradient-based optimization
 - Gradient decent/ascent
 - Newton method

Logistic Regression: Binary Classifier

Inspired by Bayes decision rule, a probabilistic formulation for classification ($y_n \in \{-1, 1\}$)¹:

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

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

or equivalently,

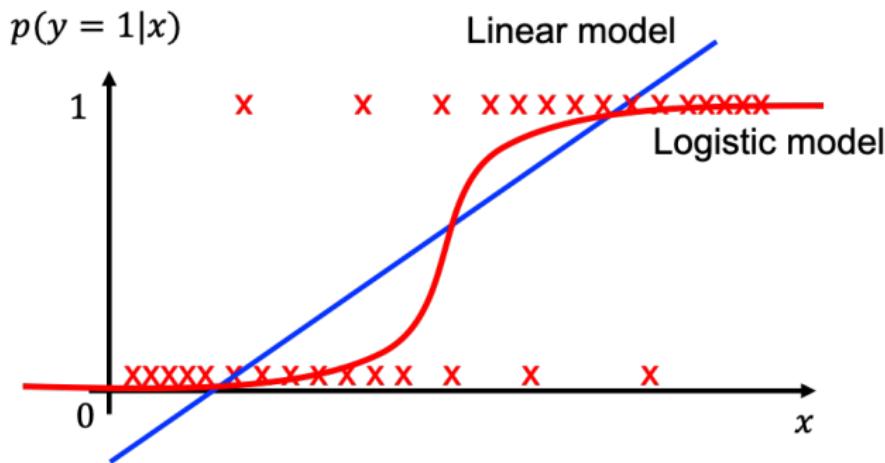
$$p(y_n \mid \mathbf{x}_n, \mathbf{w}) = \frac{1}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}$$

¹ \mathbf{x}_n can be replaced with $\phi(\mathbf{x}_n)$.

Linear vs. Logistic

- ▶ Linear regression ($y = w^\top x + b$) is easy estimation (closed form solution) but not suitable: (i) the target function (probability) is bounded $[0, 1]$, (ii) linear regression is not robust to outliers, ...
- ▶ **Logistic regression** (logistic or sigmoid function $\sigma(\cdot)$)

$$p(y = 1 | x) = \sigma(\xi(x)) = \frac{1}{1 + \exp(-\xi(x))} \in [0, 1]$$

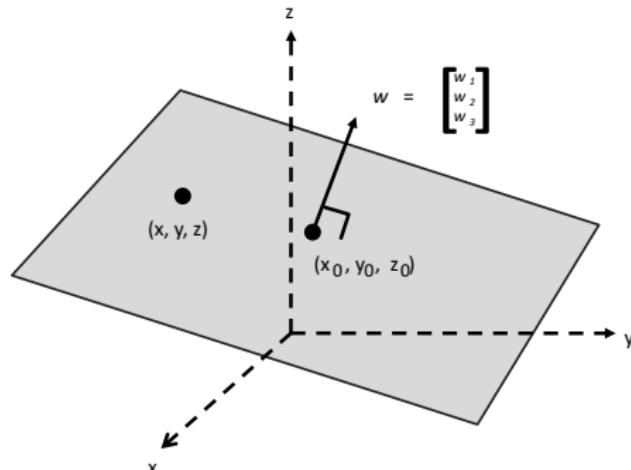


Logistic Decision Boundary

Note that we classify \mathbf{x}_n based on the sign of $-\mathbf{w}^\top \mathbf{x}_n$.

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

The decision boundary² of $\mathbf{w}^\top \mathbf{x}$:



²We assume that the last element of \mathbf{x} is 1

Property of Logistic Function

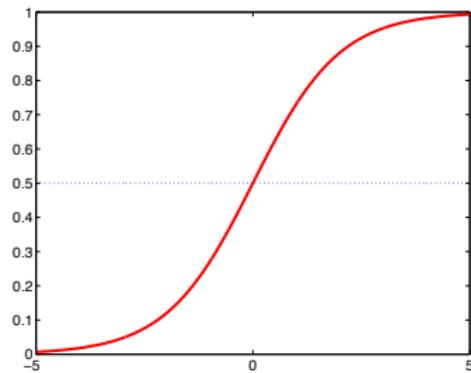


Figure: $\sigma(t) = \frac{1}{1+e^{-t}}$

- ▶ $\sigma(t) \rightarrow 0$ as $t \rightarrow -\infty$
- ▶ $\sigma(t) \rightarrow 1$ as $t \rightarrow \infty$
- ▶ $\sigma(-t) = 1 - \sigma(t)$
- ▶ $\frac{d}{dt}\sigma(t) = \sigma(t)\sigma(-t) = \sigma(t)(1 - \sigma(t))$

Optimization for Logistic Regression

Maximum likelihood estimate (to learn Bayes decision rule)

- ▶ Assume each data points are i.i.d., i.e.,

$$\log p(\mathcal{D} \mid \mathbf{w}) = \sum_{n=1}^N \log p(y_n \mid x_n, \mathbf{w})$$

- ▶ Recalling $p(y_n \mid x_n, \mathbf{w}) = \frac{1}{1+\exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}$ in logistic regression, MLE is

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

Outline

1 A systemic formulation: logistic regression

- The optimality of Bayes decision rule
- Logistic regression and probabilistic model

2 A systemic solver: gradient-based optimization

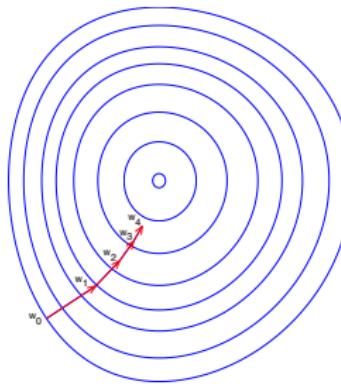
- Gradient decent/ascent
- Newton method

How to Optimize?

Logistic regression aims at finding

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

- ▶ Unfortunately, no analytic solution³ for \mathbf{w} in general
- ▶ Gradient descent



³https://en.wikipedia.org/wiki/Closed-form_expression

How to Optimize?

Logistic regression aims at finding

$$\arg \min_{\mathbf{w}} \underbrace{\sum_{n=1}^N \log (1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))}_{=: \mathcal{L}(\mathbf{w}) \text{ or } \mathcal{L}}$$

- ▶ Gradient

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \frac{-y_n \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)} \mathbf{x}_n$$

- ▶ Gradient descent:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \alpha \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_t}$$

Alternative form: $y_n \in \{0, 1\}$

If $y_n \in \{0, 1\}$, then the likelihood becomes

$$\prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_n)^{y_n} \left(1 - \sigma(\mathbf{w}^\top \mathbf{x}_n)\right)^{(1-y_n)}$$

Then the optimization goal is

$$\arg \max_{\mathbf{w}} \sum_{n=1}^N y_n \log \sigma(\mathbf{w}^\top \mathbf{x}_n) + (1 - y_n) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_n))$$

Logistic Regression: Gradient Ascent ($y_n \in \{0, 1\}$)

The gradient ascent learning has the form

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \alpha \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right).$$

Calculate the gradient

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \sum_{n=1}^N [y_n - \sigma(\mathbf{w}^\top \mathbf{x}_n)] \mathbf{x}_n$$

Hence, the gradient ascent update rule for \mathbf{w} is

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} + \alpha \sum_{n=1}^N \left[y_n - \sigma \left((\mathbf{w}^{\text{old}})^\top \mathbf{x}_n \right) \right] \mathbf{x}_n.$$

Detailed Calculation of Gradient

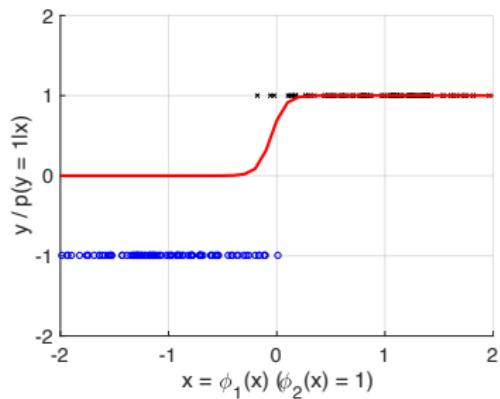
Let $\hat{y}_n = \sigma(\mathbf{w}^\top \mathbf{x}_n)$. Recall the log-likelihood

$$\mathcal{L} = \sum_{n=1}^N \left[y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n) \right].$$

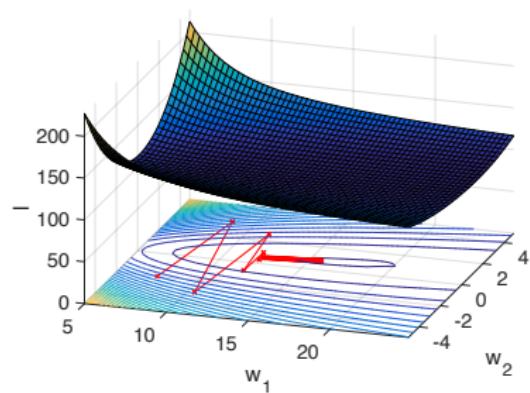
Calculate

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{w}} &= \sum_{n=1}^N \left[y_n \frac{\hat{y}'_n}{\hat{y}_n} \mathbf{x}_n + (1 - y_n) \frac{-\hat{y}'_n}{1 - \hat{y}_n} \mathbf{x}_n \right] \\ &= \sum_{n=1}^N \left[y_n \frac{\hat{y}_n(1 - \hat{y}_n)}{\hat{y}_n} - (1 - y_n) \frac{\hat{y}_n(1 - \hat{y}_n)}{1 - \hat{y}_n} \right] \mathbf{x}_n \\ &= \sum_{n=1}^N [y_n(1 - \hat{y}_n) - (1 - y_n)\hat{y}_n] \mathbf{x}_n \quad = \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n\end{aligned}$$

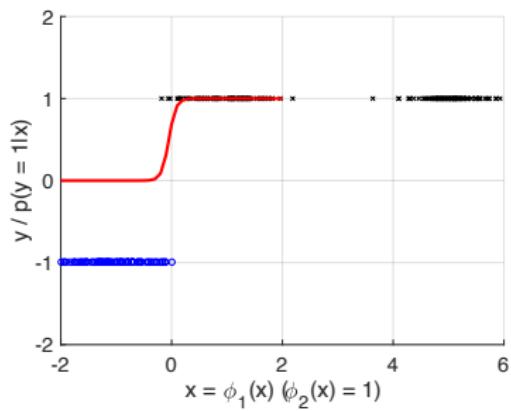
Data



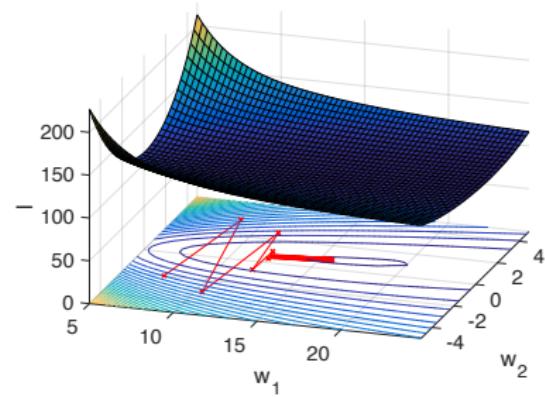
Loss



Data



Loss



Outline

1 A systemic formulation: logistic regression

- The optimality of Bayes decision rule
- Logistic regression and probabilistic model

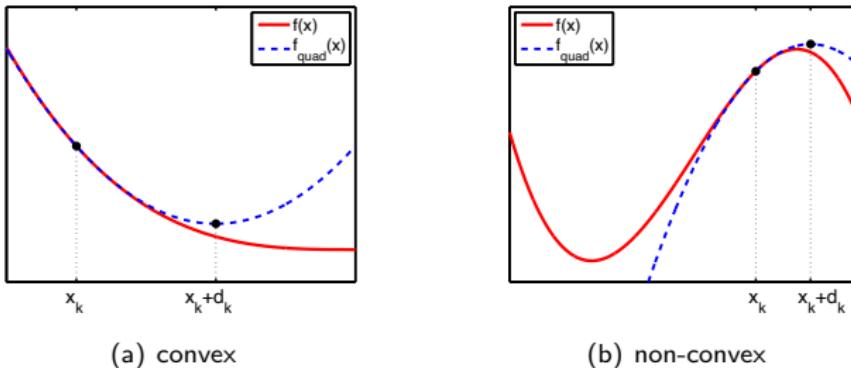
2 A systemic solver: gradient-based optimization

- Gradient decent/ascent
- Newton method

Newton's Method: Second-order Optimization

- ▶ *Iterative re-weighted least squares (IRLS)* is a popular algorithm, derived from Newton's method.
- ▶ Mathematical Preliminaries
 - ▶ Gradient
 - ▶ Gradient descent/ascent
 - ▶ *Hessian matrix*
 - ▶ *Newton's method*

Intuition behind Newton's Method



[Figure source: Murphy's book]

- ▶ $f(x)$ is our objective function (which is expensive to evaluate)
- ▶ $f_{\text{quad}}(x)$ is a function approximating $f(x)$ at x_k (relatively easy to approximate)
- ▶ We can easily obtain the minimum/maximum from $f_{\text{quad}}(x)$.

Hessian Matrix

- If $f(\mathbf{x})$ is twice differentiable, the Hessian matrix \mathbf{H} is defined as the symmetric matrix with the (i,j) element $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$

$$\begin{aligned}\mathbf{H} &= \nabla^2 f(\mathbf{x}) \\ &= \left[\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \right] \\ &= \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_D} \\ \vdots & & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_D \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_D \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_D^2} \end{bmatrix} \\ &= \frac{\partial}{\partial \mathbf{x}} \left[\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right]^\top = \frac{\partial}{\partial \mathbf{x}} [\nabla f(\mathbf{x})]^\top\end{aligned}$$

Newton's Method

The basic idea of Newton's method is to optimize the *quadratic approximation* of the objective function $\mathcal{J}(\mathbf{w})$ around the current point $\mathbf{w}^{(k)}$.

The second-order Taylor series expansion of $\mathcal{J}(\mathbf{w})$ at the current point $\mathbf{w}^{(k)}$ gives

$$\begin{aligned}\mathcal{J}_2(\mathbf{w}) &= \mathcal{J}(\mathbf{w}^{(k)}) + \left[\nabla \mathcal{J}(\mathbf{w}^{(k)}) \right]^\top (\mathbf{w} - \mathbf{w}^{(k)}) \\ &\quad \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(k)})^\top \nabla^2 \mathcal{J}(\mathbf{w}^{(k)}) (\mathbf{w} - \mathbf{w}^{(k)}),\end{aligned}$$

where $\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})$ is the Hessian of $\mathcal{J}(\mathbf{w})$ evaluated at $\mathbf{w} = \mathbf{w}^{(k)}$.

Newton's Method

The basic idea of Newton's method is to optimize the *quadratic approximation* of the objective function $\mathcal{J}(\mathbf{w})$ around the current point $\mathbf{w}^{(k)}$.

The second-order Taylor series expansion of $\mathcal{J}(\mathbf{w})$ at the current point $\mathbf{w}^{(k)}$ gives

$$\begin{aligned}\mathcal{J}_2(\mathbf{w}) &= \mathcal{J}(\mathbf{w}^{(k)}) + [\nabla \mathcal{J}(\mathbf{w}^{(k)})]^\top (\mathbf{w} - \mathbf{w}^{(k)}) \\ &\quad \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(k)})^\top \nabla^2 \mathcal{J}(\mathbf{w}^{(k)}) (\mathbf{w} - \mathbf{w}^{(k)}),\end{aligned}$$

where $\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})$ is the Hessian of $\mathcal{J}(\mathbf{w})$ evaluated at $\mathbf{w} = \mathbf{w}^{(k)}$. Differentiate $\mathcal{J}_2(\mathbf{w})$ w.r.t. \mathbf{w} and set it equal 0, which leads to

$$\nabla \mathcal{J}(\mathbf{w}^{(k)}) + \nabla^2 \mathcal{J}(\mathbf{w}^{(k)}) \mathbf{w} - \nabla^2 \mathcal{J}(\mathbf{w}^{(k)}) \mathbf{w}^{(k)} = 0$$

Thus we have

$$\mathbf{w} = \mathbf{w}^{(k)} - [\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})]^{-1} \nabla \mathcal{J}(\mathbf{w}^{(k)}).$$

Hence the Newton's method is of the form

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - [\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})]^{-1} \nabla \mathcal{J}(\mathbf{w}^{(k)}).$$

Remark

We don't need learning rate!

Remark

The Hessian $\nabla^2 \mathcal{J}(\mathbf{w}^{(k)})$ should be positive definite for all k if $\mathcal{J}(\mathbf{w}^{(k)})$ is convex.

Taylor Series

The Taylor series represents a function f as an infinite sum of terms.

Definition (Taylor Series)

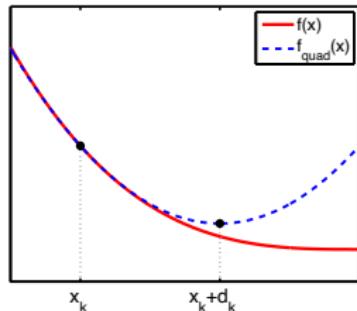
For a smooth function $f \in \mathcal{C}^\infty$, $f : \mathbb{R} \rightarrow \mathbb{R}$, the *Taylor series* of f at x_0 is defined as

$$T_\infty(x) := \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k,$$

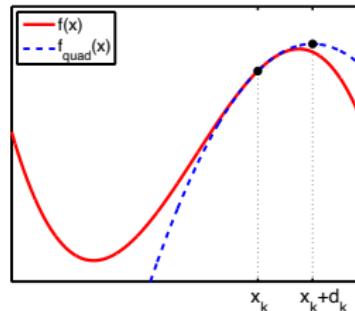
where $f^{(k)}(x_0)$ is the k th derivative of f at x_0 .

Let $T_k(x)$ is the summation up to k -th term.

Illustration of Newton's Method



(a) convex



(b) non-convex

[Figure source: Murphy's book]

Newton method requires the second order derivatives + matrix inversion!

Detailed Calculation of Hessian

Calculate the Hessian for logistic regression:

$$\begin{aligned}\nabla^2 \mathcal{L} &= \frac{\partial}{\partial \mathbf{w}} [\nabla \mathcal{L}]^\top \\ &= \frac{\partial}{\partial \mathbf{w}} \left[\sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n^\top \right] \\ &= - \sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^\top\end{aligned}$$

We set the objective function $\mathcal{J}(\mathbf{w})$ as the negative log-likelihood:

$$\mathcal{J}(\mathbf{w}) = -\mathcal{L}(\mathbf{w}) = - \sum_{n=1}^N \left[y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n) \right].$$

Thus, the gradient and the Hessian are:

$$\nabla \mathcal{J}(\mathbf{w}) = - \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n,$$

$$\nabla^2 \mathcal{J}(\mathbf{w}) = \sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x}_n \mathbf{x}_n^\top.$$

One can easily see that the Hessian is positive definite since $\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$ is positive definite and $\hat{y}_n (1 - \hat{y}_n) \geq 0, \forall n$.

Logistic Regression: IRLS

- ▶ Newton's update has the form

$$\Delta \mathbf{w} = - \underbrace{\left[\sum_{n=1}^N \hat{y}_n(1-\hat{y}_n) \mathbf{x}_n \mathbf{x}_n^\top \right]^{-1}}_{\text{inverse of Hessian}} \underbrace{\left[- \sum_{n=1}^N (y_n - \hat{y}_n) \mathbf{x}_n \right]}_{\text{gradient}}$$

- ▶ Newton's update reduces to iterative re-weighted least squares (IRLS):

$$\Delta \mathbf{w} = (\mathbf{X} \mathbf{S} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{S} \mathbf{b},$$

where

$$\mathbf{S} = \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{y}_N(1-\hat{y}_N) \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \frac{y_1 - \hat{y}_1}{\hat{y}_1(1-\hat{y}_1)} \\ \vdots \\ \frac{y_N - \hat{y}_N}{\hat{y}_N(1-\hat{y}_N)} \end{bmatrix}$$

Alternatively we write

$$\mathbf{w}_{k+1} = \mathbf{w}_k + (\mathbf{X}\mathbf{S}_k\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{S}_k\mathbf{b}_k$$

where \mathbf{S}_k and \mathbf{b}_k are computed using \mathbf{w}_k .

That is,

$$\begin{aligned}\mathbf{w}_{k+1} &= \mathbf{w}_k + (\mathbf{X}\mathbf{S}_k\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{S}_k\mathbf{b}_k \\ &= (\mathbf{X}\mathbf{S}_k\mathbf{X}^\top)^{-1} \left[(\mathbf{X}\mathbf{S}_k\mathbf{X}^\top)\mathbf{w}_k + \mathbf{X}\mathbf{S}_k\mathbf{b}_k \right] \\ &= (\mathbf{X}\mathbf{S}_k\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{S}_k \left[\mathbf{X}^\top\mathbf{w}_k + \mathbf{b}_k \right].\end{aligned}$$

Algorithm Outline

Algorithm 1 IRLS

Input: $\{(\mathbf{x}_n, y_n) \mid n = 1, \dots, N\}$

- 1: Initialize $\mathbf{w} = \mathbf{0}$ and $w_0 = \log(\bar{y}/(1 - \bar{y}))$
 - 2: **repeat**
 - 3: **for** $n = 1, \dots, N$ **do**
 - 4: Compute $\eta_n = \mathbf{w}^\top \mathbf{x}_n + w_0$
 - 5: Compute $\hat{y}_n = \sigma(\eta_n)$
 - 6: Compute $s_n = \hat{y}_n(1 - \hat{y}_n)$
 - 7: Compute $z_n = \eta_n + \frac{\eta_n - \hat{y}_n}{s_n}$
 - 8: **end for**
 - 9: Construct $\mathbf{S} = \text{diag}(s_{1:N})$
 - 10: Update $\mathbf{w} = (\mathbf{X}\mathbf{S}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{S}\mathbf{z}$
 - 11: **until** convergence
 - 12: **return** \mathbf{w}
-

Summary

Classification

- ▶ Issues when using linear regression for classification
- ▶ A systemic formulation: logistic regression
 - ▶ The optimality of Bayes decision rule
 - ▶ Logistic regression and probabilistic model
- ▶ A systemic solver: gradient-based optimization
 - ▶ Gradient decent/ascent
 - ▶ Newton method (IRLS for logistic regression)

Multiclass Logistic Regression

- ▶ Model input-output by a softmax transformation of linear functions of input features:

$$p(y_n = k \mid x_n) = \frac{\exp(\mathbf{w}_k^\top \mathbf{x}_n)}{\sum_{j=1}^K \exp(\mathbf{w}_j^\top \mathbf{x}_n)}$$

- ▶ The log-likelihood is given by

$$\sum_{n=1}^N \sum_{k=1}^K y_{kn} \log p(y_n = k \mid x_n)$$

where $y_{kn} = \mathbb{1}[y_n = k]$.

Further Readings

- ▶ Textbook: Chapter 10 (Logistic Regression)
- ▶ Textbook: Chapter 8.3 (Optimization: Second-order methods)