



Pattern Recognition and Image Analysis

Dr. Manal Helal – Fall 2015
Lecture 4



Solving Posteriors (Deterministic)

Naive Bayes

Naïve Bayes and Statistical Independence

Naïve Bayes is

- still widely (and successfully) used
- often outperforming much more advanced classifiers
- appropriate in the presence of high dimensional features (curse of dimensionality)
- also called “Idiot’s Bayes”

Naïve Bayes and Statistical Independence (cont.)

For the class dependent pdf we can do the following factorization:

$$\begin{aligned} p(\mathbf{x}|y) &= p(x_1, x_2, \dots, x_d|y) \\ &= p(x_1|y)p(x_2, x_3, \dots, x_d|y, x_1) \\ &= p(x_1|y)p(x_2|y, x_1)p(x_3, x_4, \dots, x_d|y, x_1, x_2) \\ &= p(x_1|y) \prod_{i=2}^d p(x_i|y, x_1, \dots, x_{i-1}) \end{aligned}$$

Naïve Bayes and Statistical Independence (cont.)

- The Naïve Bayes classifier makes a very strong – so to call naïve – independency assumption.
- All d components of the feature vector \mathbf{x} are assumed to be mutually independent.
- This independency assumption implies:

$$p(\mathbf{x}|y) = \prod_{i=1}^d p(x_i|y)$$

Naïve Bayes and Statistical Independence (cont.)

The decision rule of naïve Bayes reads as follows:

$$\begin{aligned}y^* &= \operatorname{argmax}_y p(y|\mathbf{x}) \\&= \operatorname{argmax}_y p(y)p(\mathbf{x}|y) \\&= \operatorname{argmax}_y p(y) \prod_{i=1}^d p(x_i|y)\end{aligned}$$

An Example: Naïve Bayes and Gaussians

Example

Assume the 100–dimensional feature vector $\mathbf{x} \in \mathbb{R}^{100}$ belonging to class y is normally distributed and all components are *mutually dependent*:

$$\begin{aligned}\boldsymbol{\mu}_y &\in \mathbb{R}^{100} \\ \boldsymbol{\Sigma} &= \boldsymbol{\Sigma}^T \in \mathbb{R}^{100 \times 100}\end{aligned}$$

The total number of parameters to be estimated for each class is

$$100 + 100 \cdot (100 + 1)/2 = 5150.$$

An Example: Naïve Bayes and Gaussians (cont.)

Example cont.

Assume the 100–dimensional feature vector $\mathbf{x} \in \mathbb{R}^{100}$ belonging to class y is normally distributed and all components are *mutually independent*.

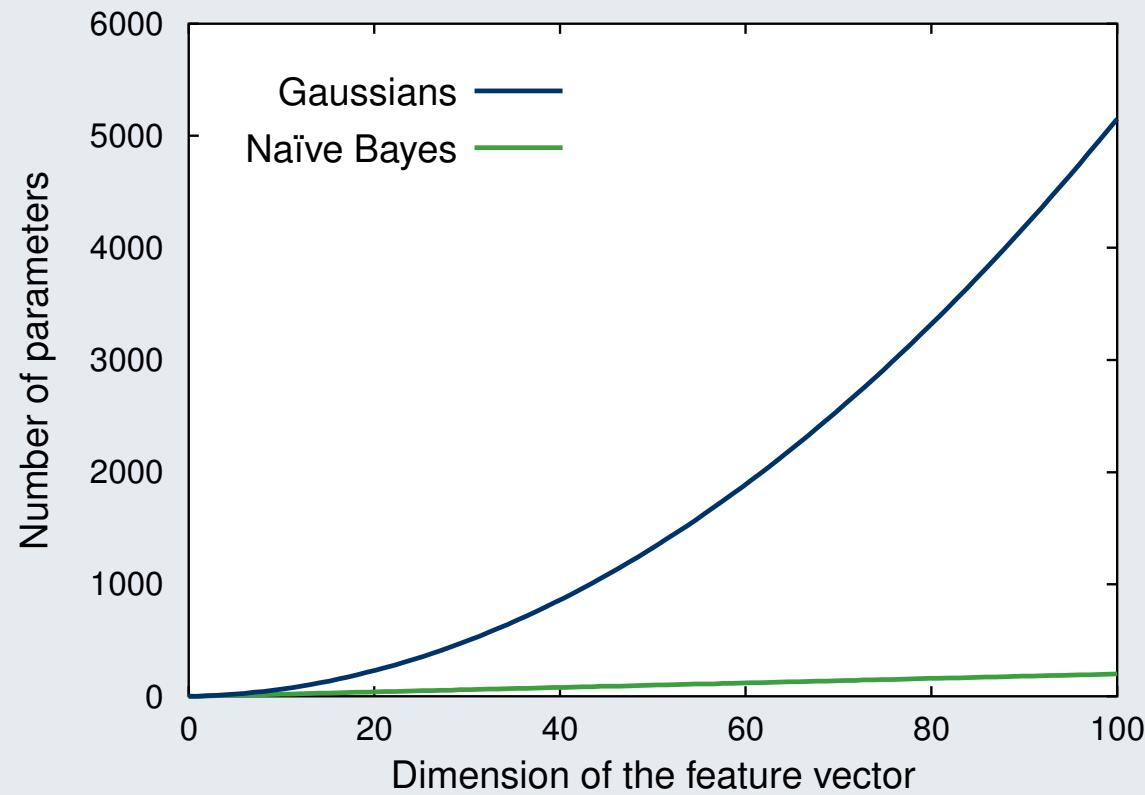
$$p(\mathbf{x}|y) = \prod_{i=1}^{100} p(x_i|y) = \prod_{i=1}^{100} \mathcal{N}(x_i; \mu_i, \sigma_i^2).$$

For each component $i = \{1, 2, 3, \dots, 100\}$ we have to estimate mean $\mu_i \in \mathbb{R}$ and variance $\sigma_i^2 \in \mathbb{R}$. The total number of parameters to be estimated for each class is

$$100 + 100 = 200.$$

An Example: Naïve Bayes and Gaussians (cont.)

Example cont.



An Example: Naïve Bayes and Gaussians (cont.)

Example cont.

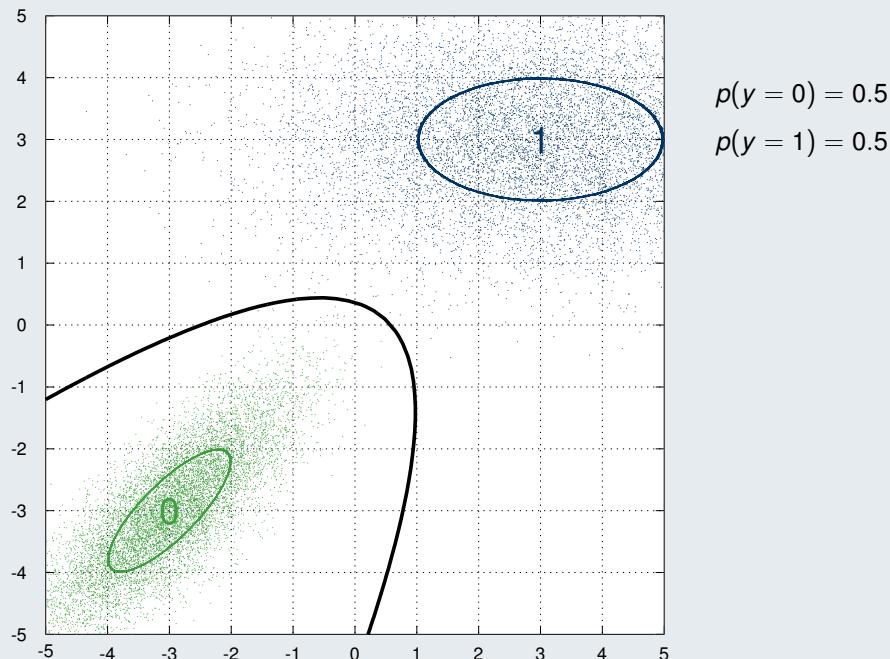


Fig.: Quadratic decision boundary that considers statistical dependency

An Example: Naïve Bayes and Gaussians (cont.)

Example cont.

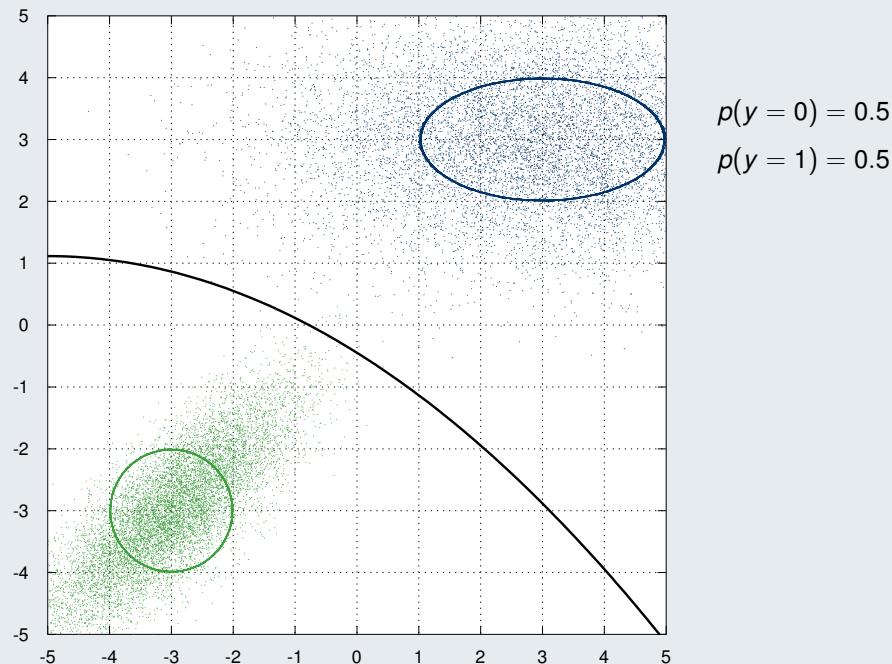


Fig.: Quadratic decision boundary assuming independency of x_1 and x_2

Naïve Bayes

Let us consider the logit transform

$$\begin{aligned}
 \log \frac{p(y=0|\mathbf{x})}{p(y=1|\mathbf{x})} &= \log \frac{p(y=0)p(\mathbf{x}|y=0)}{p(y=1)p(\mathbf{x}|y=1)} \\
 &= \log \frac{p(y=0)}{p(y=1)} + \log \frac{p(\mathbf{x}|y=0)}{p(\mathbf{x}|y=1)} \\
 &= \log \frac{p(y=0)}{p(y=1)} + \log \frac{\prod_{i=1}^d p(x_i|y=0)}{\prod_{i=1}^d p(x_i|y=1)} \\
 &= \underbrace{\alpha_0 + \sum_{i=1}^d \alpha_{0,i}(x_i)}_{\text{generalized additive model}}
 \end{aligned}$$

Naïve Bayes (cont.)

Is there anything between Bayes and Naïve Bayes?

Naïve Bayes (cont.)

There are multiple techniques to beat the curse of dimensionality, for example:

- Reduction of the parameter space
 - Introduction of independency assumptions
(from complete dependency to mutual independency)
 - Parameter tying
- Reduction of the dimension of the feature vectors

Naïve Bayes (cont.)

First order dependency

$$\begin{aligned} p(\mathbf{x}|y) &= p(x_1, x_2, \dots, x_d|y) \\ &= p(x_1|y)p(x_2, x_3, \dots, x_d|y, x_1) \\ &= p(x_1|y)p(x_2|y, x_1)p(x_3, x_4, \dots, x_d|y, x_1, x_2) \\ &= p(x_1|y) \prod_{i=2}^d p(x_i|y, x_{i-1}) \end{aligned}$$

Naïve Bayes (cont.)

Example

First order dependency in a Gaussian random vector can be identified through the covariance matrix Σ . It has the following structure:

$$\Sigma = \begin{pmatrix} \sigma_{1,1} & \sigma_{2,1} & 0 & 0 & \cdots & 0 & 0 \\ \sigma_{2,1} & \sigma_{2,2} & \sigma_{3,2} & 0 & \cdots & 0 & 0 \\ 0 & \sigma_{3,2} & \sigma_{3,3} & \sigma_{4,3} & \cdots & 0 & 0 \\ 0 & 0 & \sigma_{4,3} & \sigma_{4,4} & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \sigma_{d,d-1} & \sigma_{d,d} \end{pmatrix}$$

Naïve Bayes (cont.)

Example

First order dependency in Gaussian random vector with tied diagonal elements, i. e. $\sigma_{i,i} = \sigma$:

$$\Sigma = \begin{pmatrix} \sigma & \sigma_{2,1} & 0 & 0 & \cdots & 0 & 0 \\ \sigma_{2,1} & \sigma & \sigma_{3,2} & 0 & \cdots & 0 & 0 \\ 0 & \sigma_{3,2} & \sigma & \sigma_{4,3} & \cdots & 0 & 0 \\ 0 & 0 & \sigma_{4,3} & \sigma & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \sigma_{d,d-1} & \sigma \end{pmatrix}$$



Maximum Likelihood Estimation

Parameters Estimation

Parameterization

- Until now, $F(\mathbf{x})$ was some arbitrary function in \mathbf{x}
Example: $F(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + \boldsymbol{\alpha}^T \mathbf{x} + \alpha_0$ with components defined by Gaussian distributions
- We can express a nonlinear $F(\mathbf{x})$ as a scalar product by lifting \mathbf{x} into a higher dimensional space:

Given

$$\mathbf{x} = (x_1, x_2)^T \in \mathbb{R}^2, \quad \mathbf{A} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}, \quad \boldsymbol{\alpha} = (\alpha_1, \alpha_2)^T, \quad \alpha_0,$$

then

$$F(\mathbf{x}) = a_{11}x_1^2 + (a_{12} + a_{21})x_1x_2 + a_{22}x_2^2 + \alpha_1x_1 + \alpha_2x_2 + \alpha_0.$$

- Rewrite $F(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}'$ with $\boldsymbol{\theta}, \mathbf{x}' \in \mathbb{R}^6$:
- $$\boldsymbol{\theta} = (a_{11}, a_{12} + a_{21}, a_{22}, \alpha_1, \alpha_2, \alpha_0)^T$$
- $$\mathbf{x}' = (x_1^2, x_1x_2, x_2^2, x_1, x_2, 1)^T$$

Parameterization (cont.)

Definition

We write the parameterized logistic function in the following:

$$g(\theta^T \mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

where θ, \mathbf{x} are the lifted parameters of the original decision function F (if it was not already linear).

Log-Likelihood Function

- Let us assume the posteriors are given by

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

where $g(\boldsymbol{\theta}^T \mathbf{x})$ is the sigmoid function parameterized in $\boldsymbol{\theta}$.

- The parameter vector $\boldsymbol{\theta}$ has to be estimated from a set S of m training samples:

$$S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), \dots, (\mathbf{x}_m, y_m)\} .$$

- Method of choice: Maximum Likelihood Estimation

Log-Likelihood Function (cont.)

Before we work on the formulas of the ML-estimator, we rewrite the posteriors using Bernoulli probability:

$$p(y|\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x})^y (1 - g(\boldsymbol{\theta}^T \mathbf{x}))^{1-y}$$

which shows the great benefit of the chosen notation for class numbers.

Log-Likelihood Function (cont.)

Now we can compute the log-likelihood function
(assuming that the training samples are mutually independent):

$$\begin{aligned}\mathcal{L}(\theta) &= \log \left(\prod_{i=1}^m p(y_i | \mathbf{x}_i) \right) \\ &= \sum_{i=1}^m \log \left(g(\theta^T \mathbf{x}_i)^{y_i} (1 - g(\theta^T \mathbf{x}_i))^{1-y_i} \right) \\ &= \sum_{i=1}^m \left(y_i \log g(\theta^T \mathbf{x}_i) + (1 - y_i) \log (1 - g(\theta^T \mathbf{x}_i)) \right)\end{aligned}$$

Log-Likelihood Function (cont.)

Simplification of the log-likelihood function:

$$\begin{aligned}
 \mathcal{L}(\theta) &= \sum_{i=1}^m \left(y_i \log g(\theta^T \mathbf{x}_i) + (1 - y_i) \log (1 - g(\theta^T \mathbf{x}_i)) \right) \\
 &= \sum_{i=1}^m \left(y_i \log \frac{e^{\theta^T \mathbf{x}_i}}{1 + e^{\theta^T \mathbf{x}_i}} + (1 - y_i) \log \frac{1}{1 + e^{\theta^T \mathbf{x}_i}} \right) \\
 &= \sum_{i=1}^m \left(y_i \theta^T \mathbf{x}_i + \log \frac{1}{1 + e^{\theta^T \mathbf{x}_i}} \right) \\
 &= \sum_{i=1}^m \left(y_i \theta^T \mathbf{x}_i + \log (1 - g(\theta^T \mathbf{x}_i)) \right)
 \end{aligned}$$

Maximization of the log-likelihood function

- The log-likelihood function is concave.
- We use the  Newton-Raphson algorithm to solve the unconstrained optimization problem:

For the $(k + 1)$ -st iteration step, we get:

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \left(\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \mathcal{L}(\boldsymbol{\theta}^{(k)}) \right)^{-1} \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}^{(k)})$$

Note: If you write the Newton-Raphson iteration in matrix form, you will end up with a weighted least squares iteration scheme.

Newton-Raphson Iteration

Taylor's Theorem:

Approximation of a k -times differentiable function $f(x)$ around a given point x_0 :

$$f(x_0+h) = f(x_0) + f'(x_0)h + \frac{f''(x_0)}{2!}h^2 + \dots + \frac{f^{(k)}(x_0)}{k!}h^k + r_k(x_0+h)h^k, \quad \lim_{h \rightarrow 0} r_k(x_0+h) = 0$$

Second order Taylor polynomial:

$$f(x_0 + h) \approx f(x_0) + f'(x_0)h + \frac{1}{2}f''(x_0)h^2$$

Newton-Raphson Iteration (cont.)

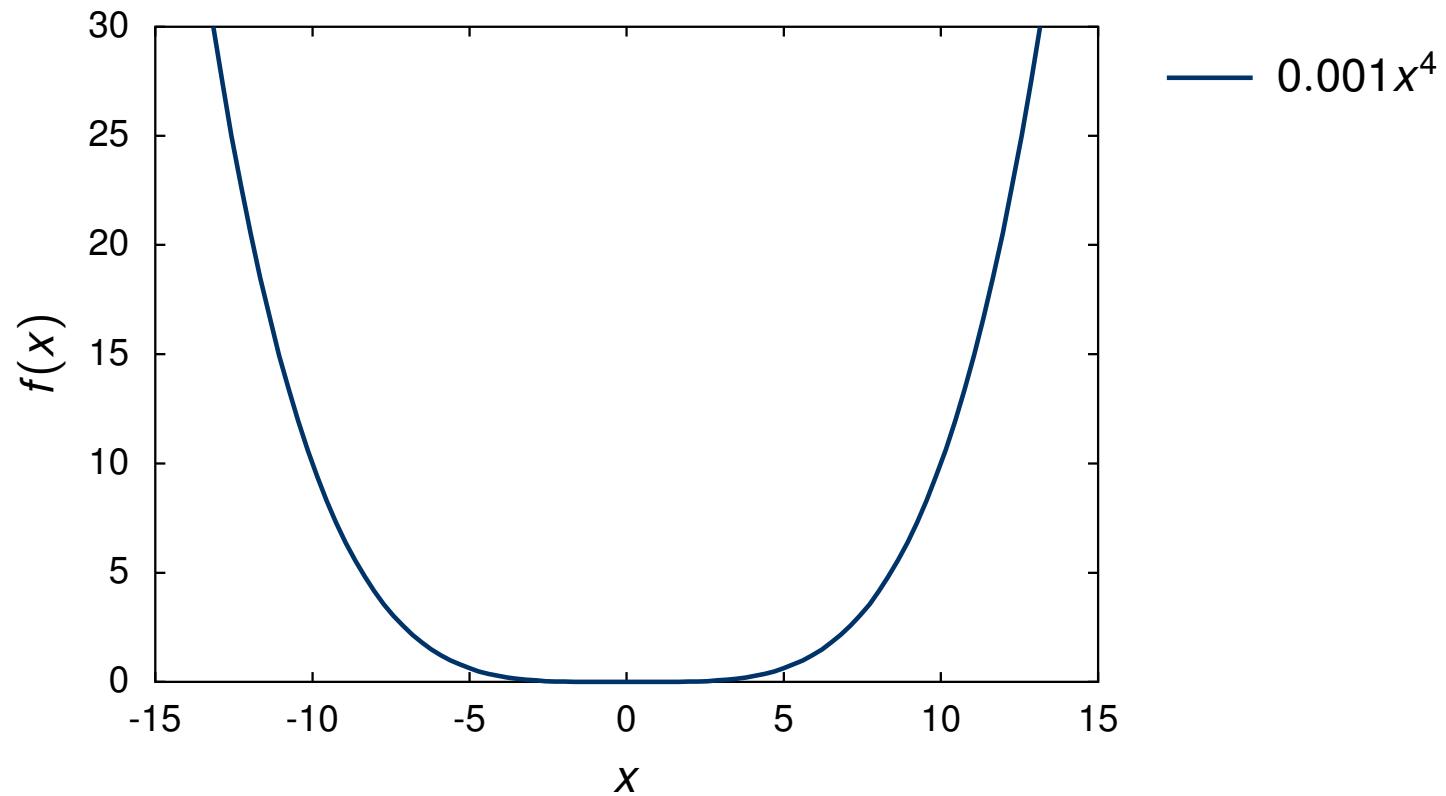
Extremum:

$$f'(x_0 + h) = f'(x_0) + f''(x_0)h \stackrel{!}{=} 0$$

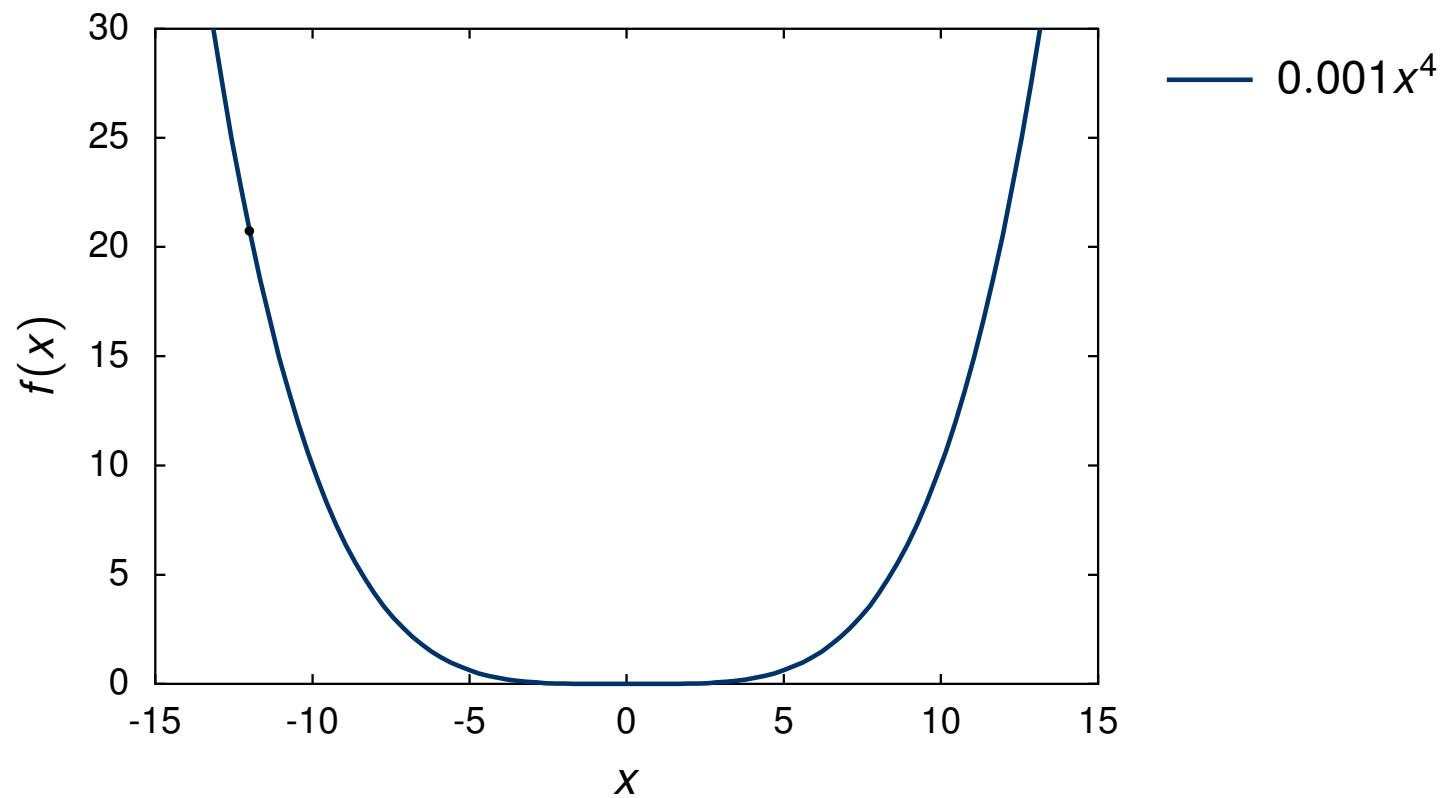
$$\hat{h} = -\frac{f'(x_0)}{f''(x_0)}$$

$$x_1 = x_0 + \hat{h} = x_0 - \frac{f'(x_0)}{f''(x_0)}$$

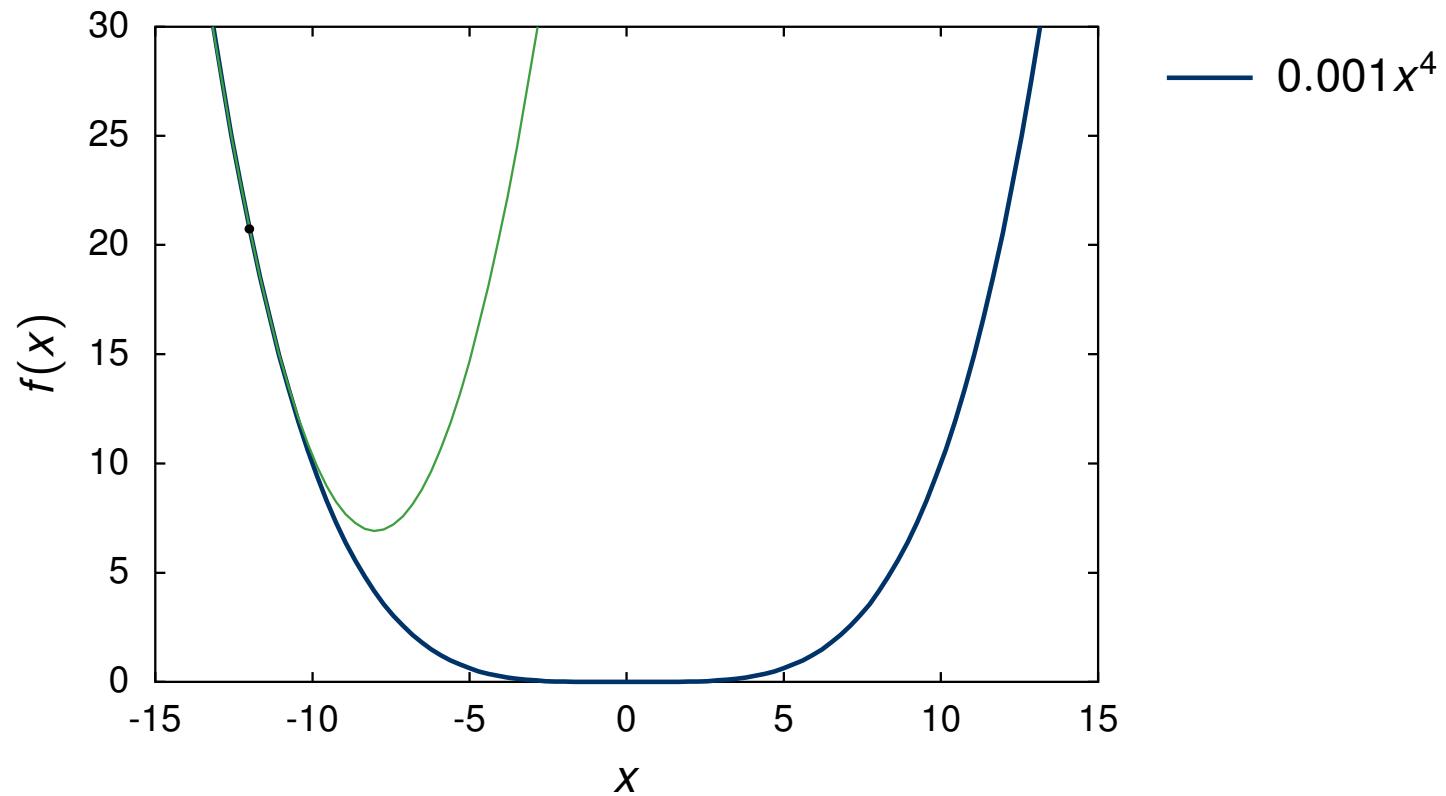
Newton-Raphson Iteration (cont.)



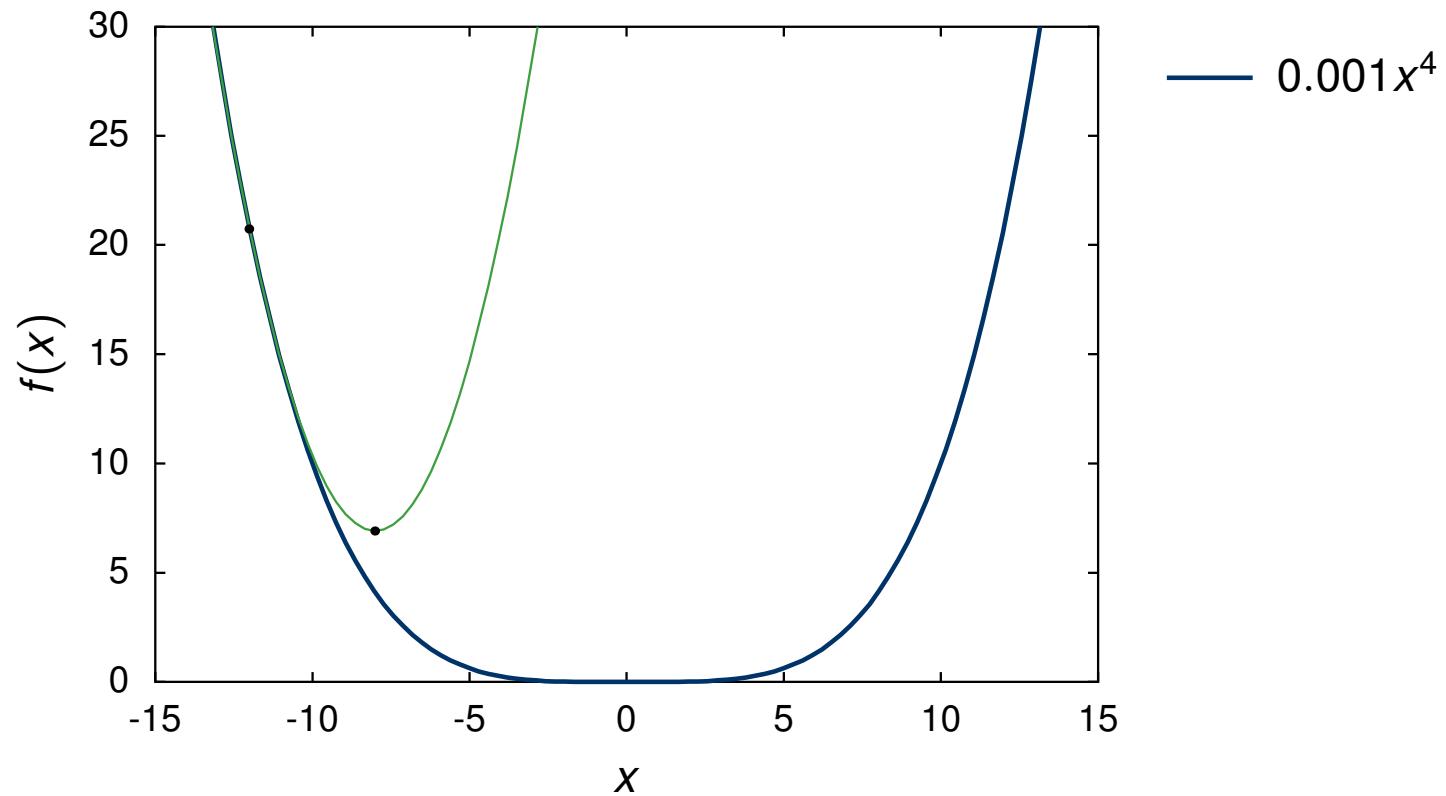
Newton-Raphson Iteration (cont.)



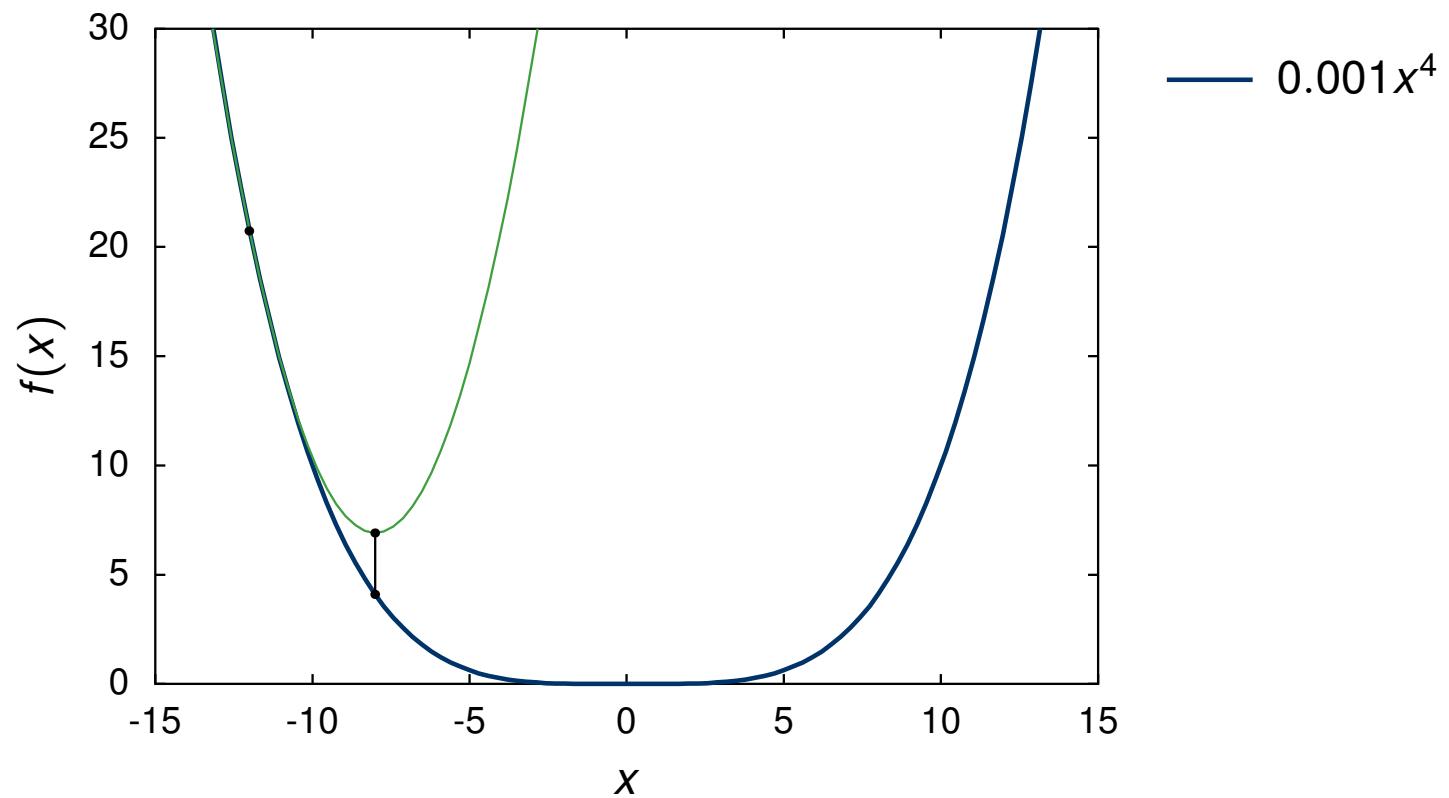
Newton-Raphson Iteration (cont.)



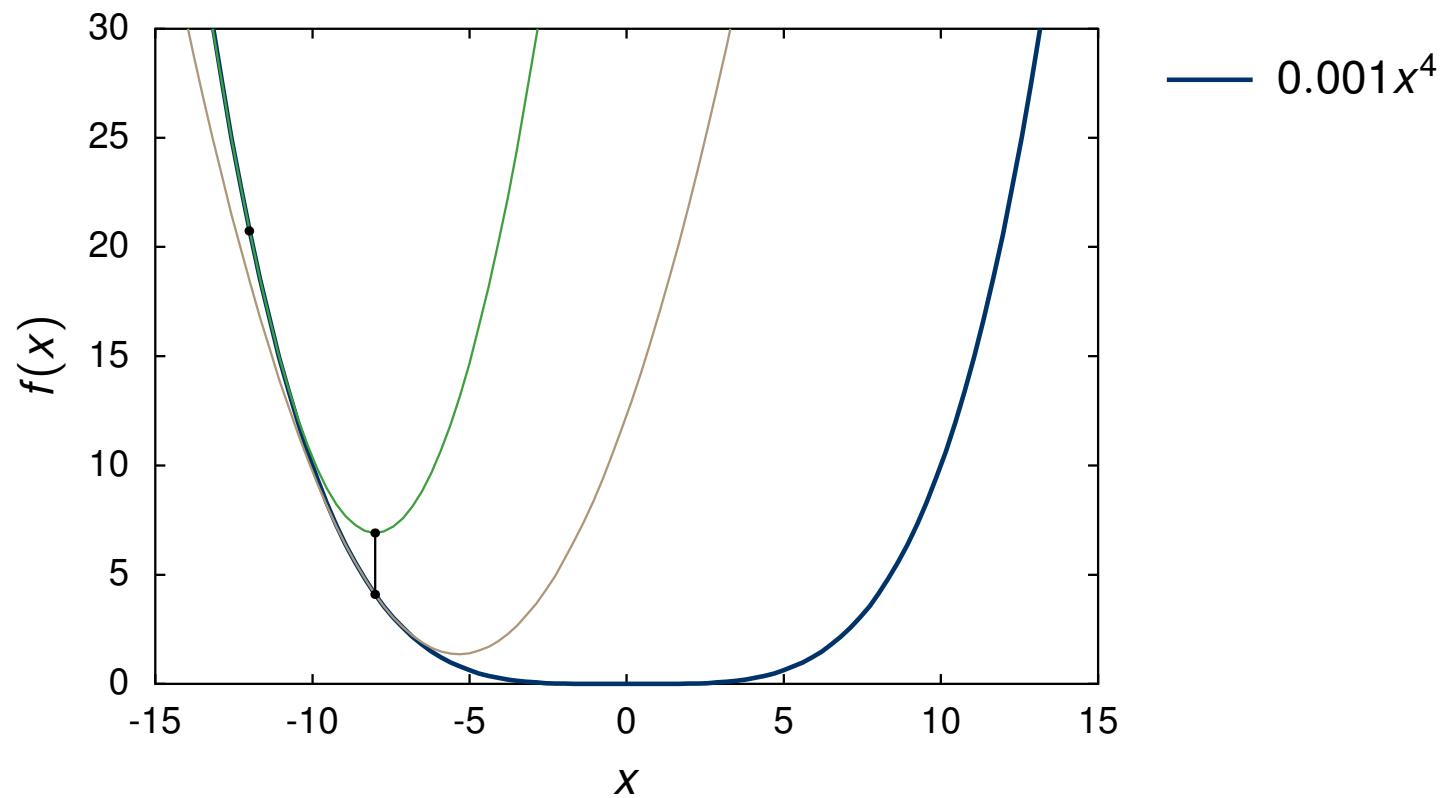
Newton-Raphson Iteration (cont.)



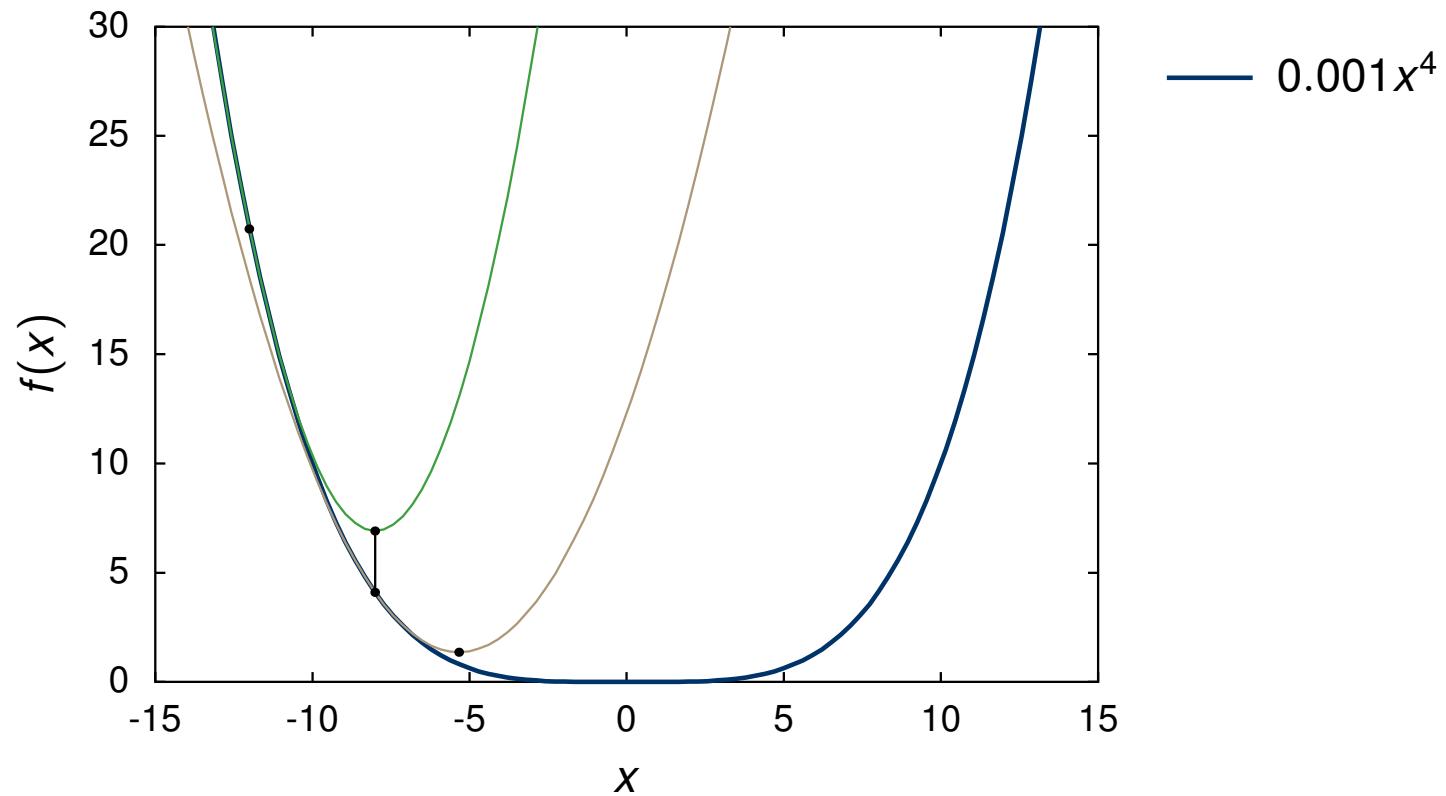
Newton-Raphson Iteration (cont.)



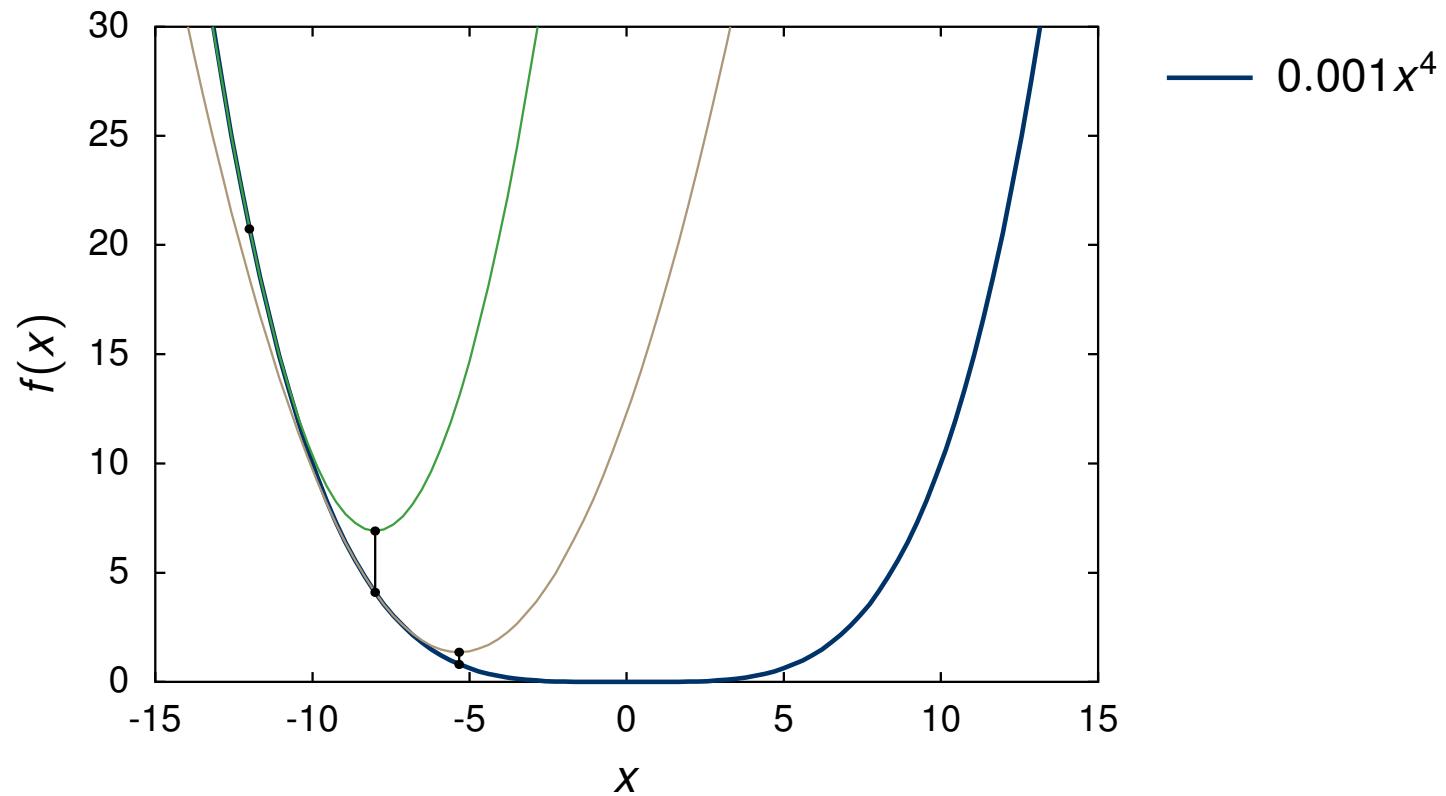
Newton-Raphson Iteration (cont.)



Newton-Raphson Iteration (cont.)



Newton-Raphson Iteration (cont.)



Gradient of the Log-Likelihood Function

The gradient:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta) &= \frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^m \left(y_i \theta^T \mathbf{x}_i + \log(1 - g(\theta^T \mathbf{x}_i)) \right) \right) \\ &= \sum_{i=1}^m \left(y_i x_{i,j} - \frac{1}{1 - g(\theta^T \mathbf{x}_i)} \frac{\partial}{\partial \theta_j} g(\theta^T \mathbf{x}_i) \right)\end{aligned}$$

Now we use the derivative of the sigmoid function and get

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta) &= \sum_{i=1}^m \left(y_i x_{i,j} - \frac{1}{1 - g(\theta^T \mathbf{x}_i)} g(\theta^T \mathbf{x}_i)(1 - g(\theta^T \mathbf{x}_i)) x_{i,j} \right) \\ &= \sum_{i=1}^m \left(y_i - g(\theta^T \mathbf{x}_i) \right) x_{i,j}\end{aligned}$$

where $x_{i,j}$ is the j -th component of the i -th training feature vector.

Gradient of the Log-Likelihood Function (cont.)

Finally, we have a quite simple gradient:

$$\frac{\partial}{\partial \theta_j} \mathcal{L}(\theta) = \sum_{i=1}^m (y_i - g(\theta^T \mathbf{x}_i)) x_{i,j}$$

where $x_{i,j}$ is the j -th component of the i -th training feature vector.

Or in vector notation:

$$\frac{\partial}{\partial \theta} \mathcal{L}(\theta) = \sum_{i=1}^m (y_i - g(\theta^T \mathbf{x}_i)) \mathbf{x}_i$$

Hessian of the Log-Likelihood Function

- The Newton-Raphson algorithm requires the Hessian matrix.
- Remember the derivative of the sigmoid function!

$$\frac{\partial^2}{\partial \theta \partial \theta^T} \mathcal{L}(\theta) = - \sum_{i=1}^m g(\theta^T \mathbf{x}_i) \left(1 - g(\theta^T \mathbf{x}_i)\right) \mathbf{x}_i \mathbf{x}_i^T$$

Example



Comprehensive Questions (cont.)

- How can a nonlinear function be written as a scalar product?
- What is the objective function for the ML-estimation of the logistic regression parameters?
- What is the difference between a gradient descent and Newton-Raphson numerical optimization scheme?
- What is the parameter update rule for the logistic regression parameters using the Newton-Raphson scheme?

Comprehensive Questions

- What is the assumption of Naïve Bayes?
- How does the assumption affect the class dependent pdf?
- What is the structure of the covariance matrix of normal-distributed classes in Naïve Bayes?
- How can Naïve Bayes be extended to first-order statistical dependencies?