

Classification and Regression III

Krzysztof Dembczyński

Intelligent Decision Support Systems Laboratory (IDSS)
Poznań University of Technology, Poland



Software Development Technologies
Master studies, second semester
Academic year 2018/19 (winter course)

Review of the previous lectures

- Mining of massive datasets.
- Classification and regression
 - ▶ What is machine learning?
 - ▶ Supervised learning: statistical decision/learning theory, loss functions, risk.
 - ▶ Learning paradigms and principles.
 - ▶ Learning algorithms: lazy learning, decision trees, generative models, linear models.

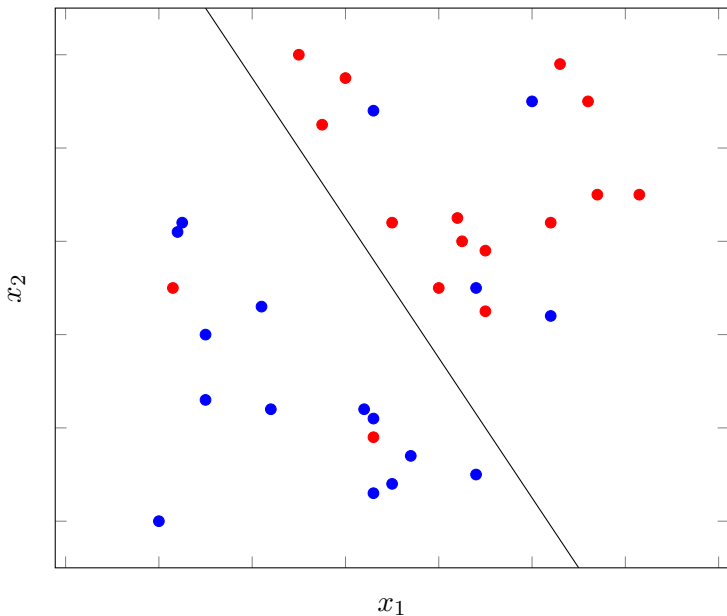
Outline

- 1 Linear Models for Classification
- 2 Summary

Outline

- 1 Linear Models for Classification
- 2 Summary

Linear models for classification



Linear models for classification

- Let the output variable be $y \in \{-1, 1\}$ and prediction function $h(\mathbf{x}) \in \{-1, 1\}$.

Linear models for classification

- Let the output variable be $y \in \{-1, 1\}$ and prediction function $h(\mathbf{x}) \in \{-1, 1\}$.

- Alternatively, one can assume $y, h(\mathbf{x}) \in \{0, 1\}$.
- Mapping $m : \{-1, 1\} \longrightarrow \{0, 1\}$ and $m^{-1} : \{0, 1\} \longrightarrow \{-1, 1\}$:

$$m(y) = \frac{y + 1}{2}, \quad m^{-1}(y) = 2y - 1.$$

Linear models for classification

- Loss is measured usually in terms of 0/1 loss which can be expressed by:

$$\ell_{0/1}(y, h(\mathbf{x})) = \mathbb{I}[yh(\mathbf{x}) \leq 0]$$

Linear models for classification

- Loss is measured usually in terms of 0/1 loss which can be expressed by:

$$\ell_{0/1}(y, h(\mathbf{x})) = \mathbb{I}[yh(\mathbf{x}) \leq 0]$$

- Solve:

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell_{0/1}(y_i, h(\mathbf{x}_i)) = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i h(\mathbf{x}_i) \leq 0]$$

Linear models for classification

- Loss is measured usually in terms of 0/1 loss which can be expressed by:

$$\ell_{0/1}(y, h(\mathbf{x})) = \mathbb{I}[yh(\mathbf{x}) \leq 0]$$

- Solve:

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell_{0/1}(y_i, h(\mathbf{x}_i)) = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i h(\mathbf{x}_i) \leq 0]$$

- **Hard** to optimize h directly.

Linear models for classification

- Usually done in two phases:

Linear models for classification

- Usually done in two phases:
 - ▶ Learn a **continuous function** $f \in \mathcal{F}$:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

Linear models for classification

- Usually done in two phases:
 - ▶ Learn a **continuous function** $f \in \mathcal{F}$:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

- ▶ **Threshold** f at 0:

$$\hat{h} = \operatorname{sgn}(\hat{f}), \quad \mathcal{H} = \{h : h = \operatorname{sgn}(f), f \in \mathcal{F}\},$$

so that

$$\mathbb{I}[y_i h(\mathbf{x}_i) \leq 0] = \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0].$$

Linear models for classification

- Usually done in two phases:
 - ▶ Learn a **continuous function** $f \in \mathcal{F}$:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

- ▶ **Threshold** f at 0:

$$\hat{h} = \text{sgn}(\hat{f}), \quad \mathcal{H} = \{h: h = \text{sgn}(f), f \in \mathcal{F}\},$$

so that

$$\mathbb{I}[y_i h(\mathbf{x}_i) \leq 0] = \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0].$$

- ▶ The quantity $yf(\mathbf{x})$ is usually referred to as **margin**.

Linear models for classification

- Solve:

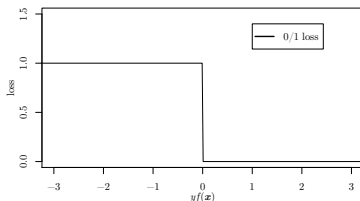
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

Linear models for classification

- Solve:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

- Still **hard** to optimize: 0/1 loss is **discontinuous** and **non-convex**.
 - ▶ e.g., when \mathcal{H} is a class of linear function, the problem known to be **NP-hard**.

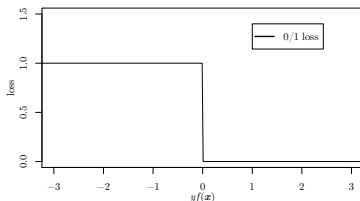


Linear models for classification

- Solve:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

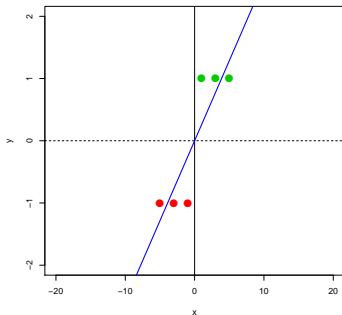
- Still **hard** to optimize: 0/1 loss is **discontinuous** and **non-convex**.
 - ▶ e.g., when \mathcal{H} is a class of linear function, the problem known to be **NP-hard**.



- **Solution:** use some **convex relaxation** of 0/1 loss.

Linear regression for classification

- We can try to use linear regression to solve binary classification problem:



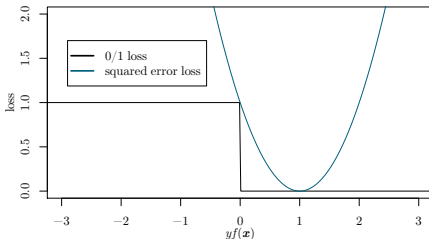
- Use $\text{sgn}(\hat{f})$ to obtain prediction.
- Minimization of squared loss leads to estimation of the conditional probability, since:

$$P(y = 1|\mathbf{x}) = \frac{\mathbb{E}(y|\mathbf{x}) + 1}{2}.$$

Linear regression for classification

- Effectively, we replace 0/1 loss by **squared error loss**:

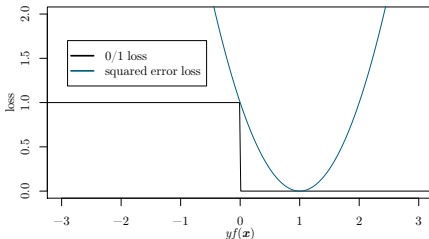
$$\ell_{\text{sq}}(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2 = (1 - yf(\mathbf{x}))^2.$$



Linear regression for classification

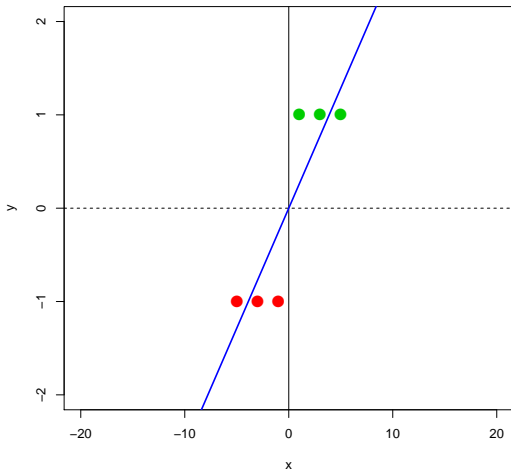
- Effectively, we replace 0/1 loss by **squared error loss**:

$$\ell_{\text{sq}}(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2 = (1 - yf(\mathbf{x}))^2.$$

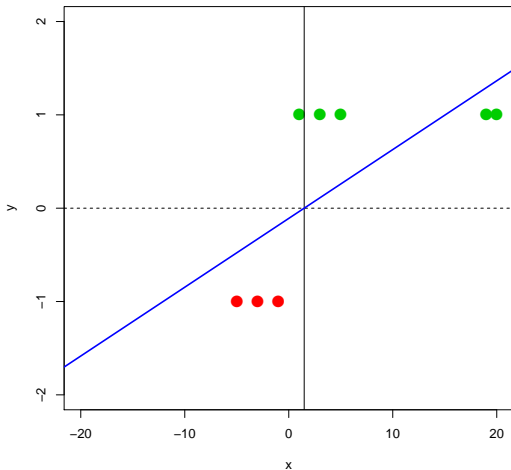


- Works nicely in practice, but has several drawbacks ...

Linear regression for classification



Linear regression for classification



- It tries to minimize the squared error of all examples, even those that are correctly classified.

Linear models for classification

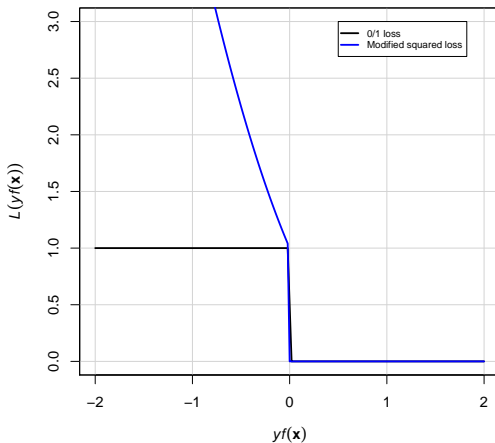
- We could consider non-zero loss only for examples with margin less or equal zero:

$$yf(\mathbf{x}) \leq 0$$

- Such a loss function could be defined as:

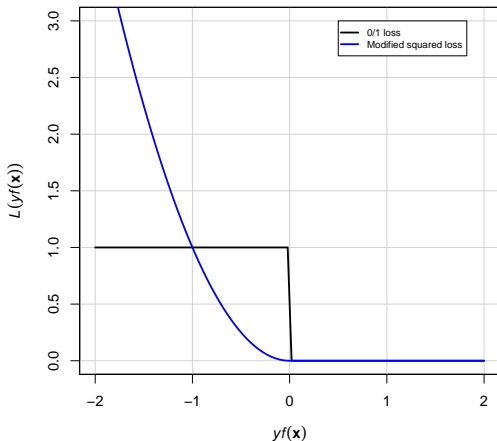
$$\ell(y, f(\mathbf{x})) = \begin{cases} 0, & \text{if } yf(\mathbf{x}) > 0 \\ (1 - yf(\mathbf{x}))^2, & \text{otherwise.} \end{cases}$$

Linear models for classification



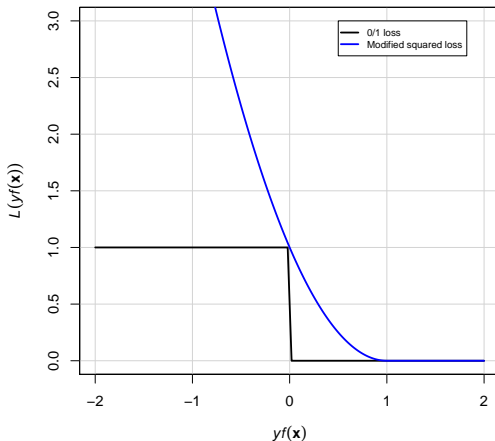
- This definition does not lead to a “nice” shape of the loss.

Linear models for classification



- A better solution: $\ell(y, f(\mathbf{x})) = (\max\{0, \epsilon - yf(\mathbf{x})\})^2$.

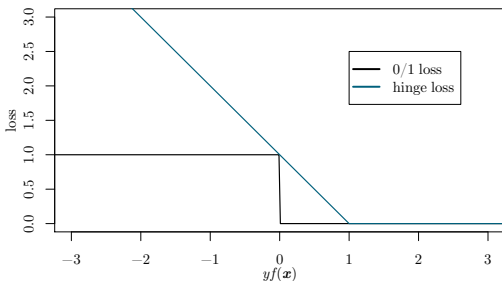
Linear models for classification



- A better solution: $\ell(y, f(\mathbf{x})) = (\max\{0, 1 - yf(\mathbf{x})\})^2$.

Linear models for classification

- Similarly, if we use absolute error instead of squared error, we get the following loss functions:
 - ▶ perceptron-like loss function: $\ell(y, f(\mathbf{x})) = \max\{0, \epsilon - yf(\mathbf{x})\}$,
 - ▶ hinge loss: $\ell(y, f(\mathbf{x})) = \max\{0, 1 - yf(\mathbf{x})\}$, used in support vector machines.



Perceptron

- Perceptron uses a linear model:

$$f(\mathbf{x}) = w_0 + \sum_{j=0}^n w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

Perceptron

- Perceptron uses a linear model:

$$f(\mathbf{x}) = w_0 + \sum_{j=0}^n w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

- We replace 0/1 loss by:

$$\begin{aligned}\ell_{\text{perc}}(y, f(\mathbf{x})) &= \max\{0, \epsilon - yf(\mathbf{x})\} \\ &= \max\{0, \epsilon - y\mathbf{w} \cdot \mathbf{x}\}\end{aligned}$$

Perceptron

- Perceptron uses a linear model:

$$f(\mathbf{x}) = w_0 + \sum_{j=0}^n w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

- We replace 0/1 loss by:

$$\begin{aligned}\ell_{\text{perc}}(y, f(\mathbf{x})) &= \max\{0, \epsilon - yf(\mathbf{x})\} \\ &= \max\{0, \epsilon - y\mathbf{w} \cdot \mathbf{x}\}\end{aligned}$$

- We solve

$$\begin{aligned}\hat{f} &= \arg \min_f \frac{1}{n} \sum_{i=1}^n \max\{0, \epsilon - y_i f(\mathbf{x}_i)\} \\ &= \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \max\{0, \epsilon - y_i \mathbf{w} \cdot \mathbf{x}_i\}\end{aligned}$$

using an incremental optimization technique (\Rightarrow the stochastic gradient descent algorithm).

Perceptron

- Learning algorithm for perceptron:

- ▶ The update in iteration t :

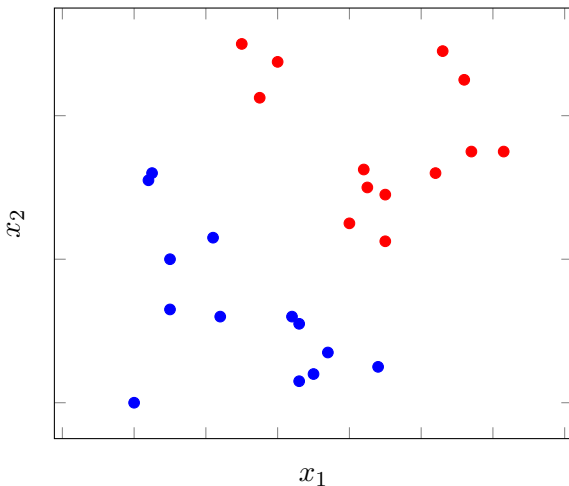
$$\begin{aligned} \mathbf{w}^t &= \mathbf{w}^{t-1} + \alpha y \mathbf{x} && \text{if } y \mathbf{w} \cdot \mathbf{x} \leq 0 \\ \mathbf{w}^t &= \mathbf{w}^{t-1} && \text{if } y \mathbf{w} \cdot \mathbf{x} > 0. \end{aligned}$$

where α is the learning rate.

- ▶ Only misclassified examples are updated.

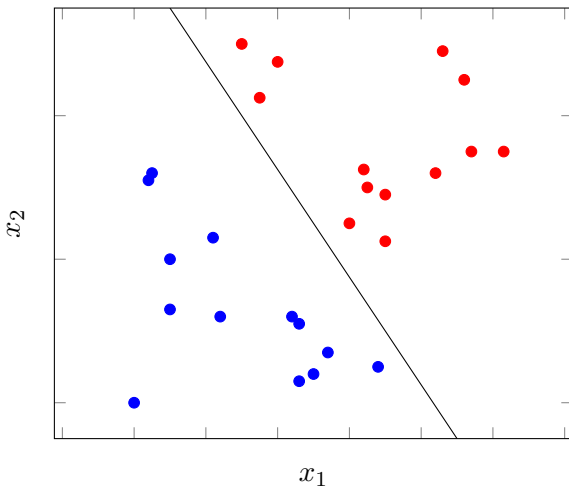
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



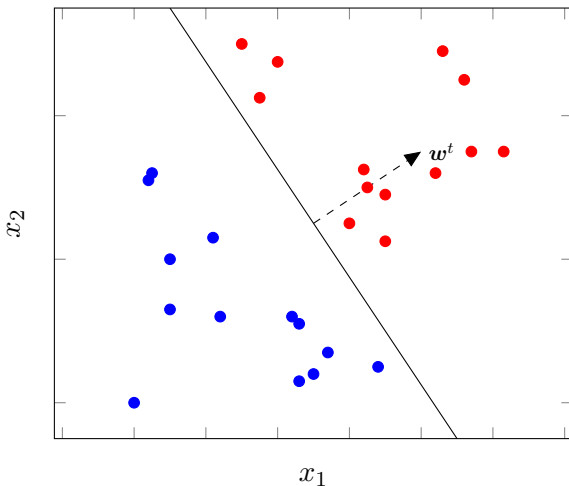
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



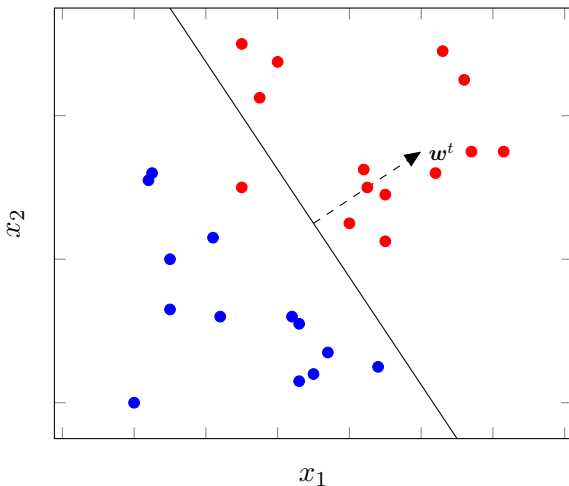
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



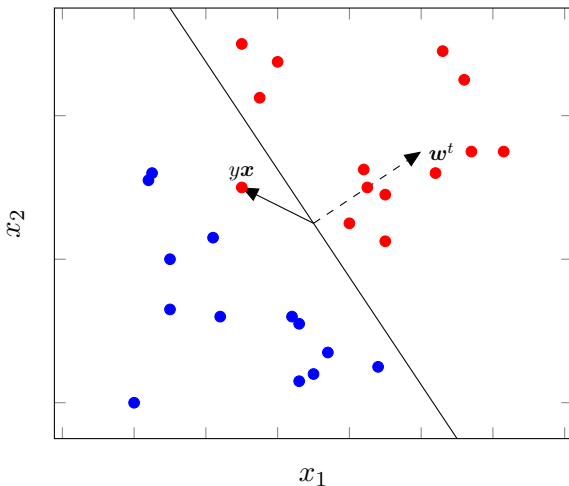
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



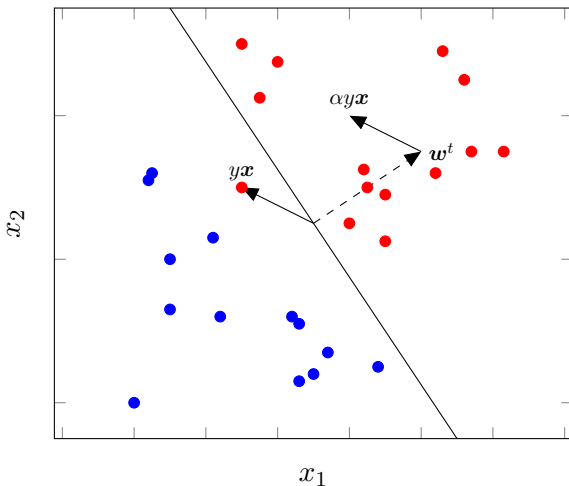
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



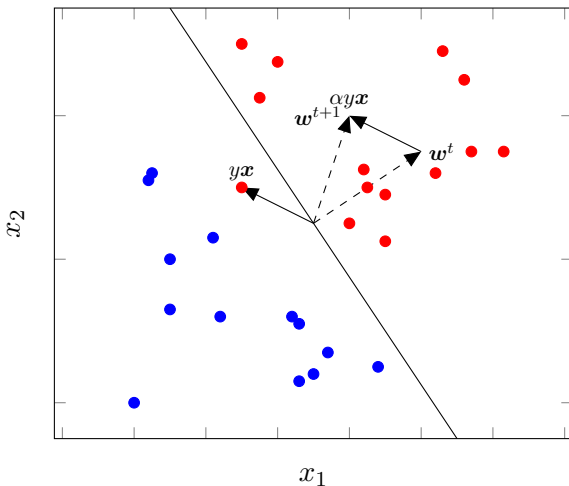
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



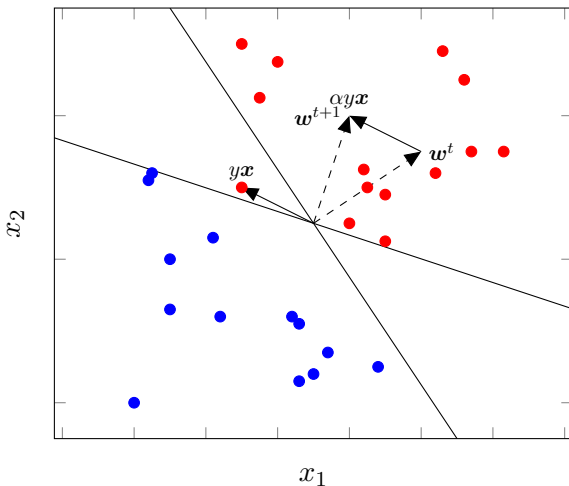
Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



Perceptron – Graphical interpretation

- The update is simply a summation or subtraction of two vectors:



Perceptron

- The update can be interpreted in terms of gradient descent:
 - ▶ For a misclassified example ($y\mathbf{w} \cdot \mathbf{x} \leq 0$), the gradient of $\ell_{\text{perc}}(y, f(\mathbf{x}))$ with respect to \mathbf{w} is given by:

$$\frac{\partial \ell_{\text{perc}}(y, f(\mathbf{x}))}{\partial \mathbf{w}} = -y\mathbf{x}.$$

- ▶ For a correctly classified example ($y\mathbf{w} \cdot \mathbf{x} > 0$), the gradient is 0.

Perceptron

- The update can be interpreted in terms of gradient descent:
 - ▶ For a misclassified example ($y\mathbf{w} \cdot \mathbf{x} \leq 0$), the gradient of $\ell_{\text{perc}}(y, f(\mathbf{x}))$ with respect to \mathbf{w} is given by:

$$\frac{\partial \ell_{\text{perc}}(y, f(\mathbf{x}))}{\partial \mathbf{w}} = -y\mathbf{x}.$$

- ▶ For a correctly classified example ($y\mathbf{w} \cdot \mathbf{x} > 0$), the gradient is 0.
 - ▶ Therefore, the update has the form:

$$\begin{aligned} \mathbf{w}^t &= \mathbf{w}^{t-1} + \alpha y\mathbf{x} && \text{if } y\mathbf{w} \cdot \mathbf{x} \leq 0 \\ \mathbf{w}^t &= \mathbf{w}^{t-1} && \text{if } y\mathbf{w} \cdot \mathbf{x} > 0. \end{aligned}$$

- Such an algorithm is usually referred to as **stochastic gradient descent**.

Stochastic gradient descent

```
Input: learning rate  $\alpha$ 
 $w = 0$ ; //(or use random values)
while (approximate minimum is obtained) {
    Randomly shuffle examples in the training set
    for  $i=1$  to  $n$  {
         $w := w - \alpha \frac{\partial \ell(y_i, f(x_i))}{\partial w}$ 
    }
}
```

Linear models with hinge loss

- Loss function:

$$\ell_{\text{hinge}}(y, f(\boldsymbol{x})) = \max\{0, 1 - yf(\boldsymbol{x})\}$$

Linear models with hinge loss

- Loss function:

$$\ell_{\text{hinge}}(y, f(\mathbf{x})) = \max\{0, 1 - yf(\mathbf{x})\}$$

- Linear model:

$$f(\mathbf{x}) = w_0 + \sum_{j=0}^n w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

Linear models with hinge loss

- Loss function:

$$\ell_{\text{hinge}}(y, f(\mathbf{x})) = \max\{0, 1 - yf(\mathbf{x})\}$$

- Linear model:

$$f(\mathbf{x}) = w_0 + \sum_{j=0}^n w_j x_j = \mathbf{w} \cdot \mathbf{x}$$

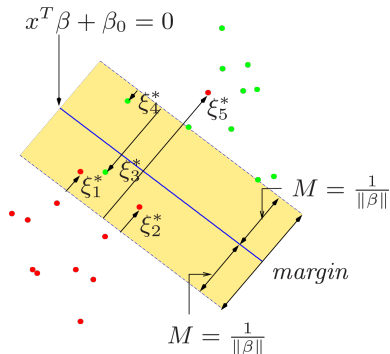
- Learning = fitting the model to the data by minimizing:

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \ell_{\text{hinge}}(y_i, f(\mathbf{x}_i)) \\ &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \max\{0, 1 - y_i f(\mathbf{x}_i)\} \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \max\{0, 1 - y_i \mathbf{w} \cdot \mathbf{x}_i\}\end{aligned}$$

Binary classification by Support Vector Machines (SVM)

Find **maximal margin classifier**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i \mathbf{w} \cdot \mathbf{x} > 1 - \xi_i \quad \forall i = 1..n \\ & \xi_i \geq 0 \quad \forall i = 1..n \end{aligned}$$



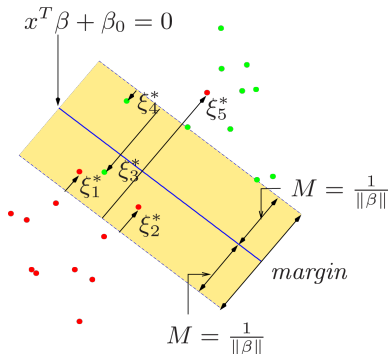
Binary classification by Support Vector Machines (SVM)

Find **maximal margin classifier**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i \mathbf{w} \cdot \mathbf{x} > 1 - \xi_i \quad \forall i = 1..n \\ & \xi_i \geq 0 \quad \forall i = 1..n \end{aligned}$$



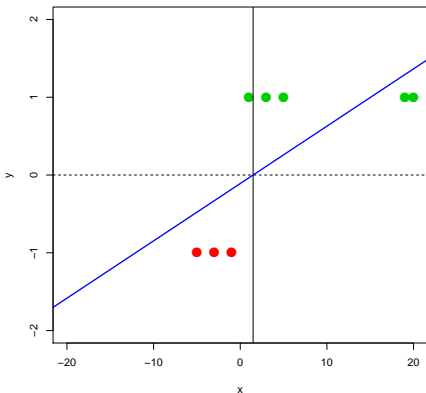
$$\begin{aligned} \min_{\mathbf{w}} \quad & \sum_{i=1}^n \max\{0, 1 - y_i \mathbf{w} \cdot \mathbf{x}_i\} \\ \text{s.t.} \quad & \|\mathbf{w}\|^2 \leq B \quad \text{for some } B. \end{aligned}$$



Logistic regression – motivation

- Another option is to use a sigmoid (or logistic) transformation of the linear function:

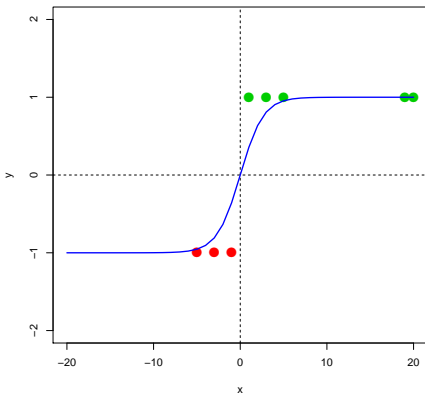
$$g(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))} \in (0, 1) \quad g(\mathbf{x}) = \frac{1 - \exp(-f(\mathbf{x}))}{1 + \exp(-f(\mathbf{x}))} \in (-1, 1)$$



Logistic regression – motivation

- Another option is to use a sigmoid (or logistic) transformation of the linear function:

$$g(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))} \in (0, 1) \quad g(\mathbf{x}) = \frac{1 - \exp(-f(\mathbf{x}))}{1 + \exp(-f(\mathbf{x}))} \in (-1, 1)$$



Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 \mid \mathbf{x})$ from $f(\mathbf{x})$.

Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 \mid \mathbf{x})$ from $f(\mathbf{x})$.
 - ▶ $\eta(\mathbf{x}) \in [0, 1]$, while $f(\mathbf{x}) \in \mathbb{R}$.

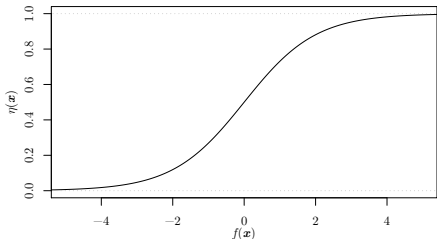
Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 \mid \mathbf{x})$ from $f(\mathbf{x})$.
 - ▶ $\eta(\mathbf{x}) \in [0, 1]$, while $f(\mathbf{x}) \in \mathbb{R}$.
 - ▶ A natural candidate for function $f(\mathbf{x}) \mapsto \eta(\mathbf{x})$: **sigmoid function**

$$\eta(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

\Updownarrow

$$f(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$$



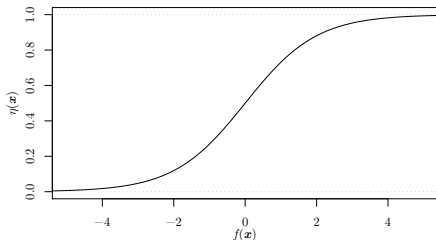
Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$ from $f(\mathbf{x})$.
 - ▶ $\eta(\mathbf{x}) \in [0, 1]$, while $f(\mathbf{x}) \in \mathbb{R}$.
 - ▶ A natural candidate for function $f(\mathbf{x}) \mapsto \eta(\mathbf{x})$: **sigmoid function**

$$\eta(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$



$$f(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$$



- Solved by the method of **Maximum Likelihood**.

$$\hat{f} = \arg \max_{f \in \mathcal{F}} P_f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$= \arg \min_{f \in \mathcal{F}} -\log P_f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i | \mathbf{x}_i)$$

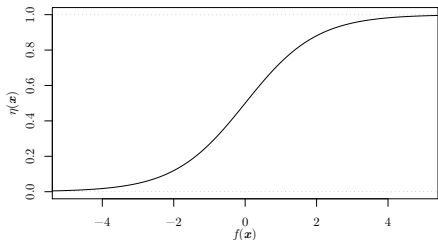
Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$ from $f(\mathbf{x})$.
 - ▶ $\eta(\mathbf{x}) \in [0, 1]$, while $f(\mathbf{x}) \in \mathbb{R}$.
 - ▶ A natural candidate for function $f(\mathbf{x}) \mapsto \eta(\mathbf{x})$: **sigmoid function**

$$\eta(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

\Updownarrow

$$f(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$$



- Solved by the method of **Maximum Likelihood**.

$$\hat{f} = \arg \max_{f \in \mathcal{F}} P_f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$= \arg \min_{f \in \mathcal{F}} -\log P_f(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i | \mathbf{x}_i)$$

Why this is correct?

Logistic regression

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i \mid \mathbf{x}_i)$$

Logistic regression

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i | \mathbf{x}_i)$$

$$\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \left(-\mathbb{I}[y_i = 1] \log \eta(\mathbf{x}_i) - \mathbb{I}[y_i = -1] \log(1 - \eta(\mathbf{x}_i)) \right)$$

Logistic regression

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i | \mathbf{x}_i)$$

$$\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \left(-\mathbb{I}[y_i = 1] \log \eta(\mathbf{x}_i) - \mathbb{I}[y_i = -1] \log(1 - \eta(\mathbf{x}_i)) \right)$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \left(\mathbb{I}[y_i = 1] \log \left(1 + e^{-f(\mathbf{x}_i)} \right) + \mathbb{I}[y_i = -1] \log \left(1 + e^{f(\mathbf{x}_i)} \right) \right)$$

$$\eta(\mathbf{x}) = \left(1 + e^{-f(\mathbf{x})} \right)^{-1}$$

$$1 - \eta(\mathbf{x}) = 1 - \frac{1}{1 + e^{-f(\mathbf{x})}} = \frac{e^{-f(\mathbf{x})}}{1 + e^{-f(\mathbf{x})}} = \left(1 + e^{f(\mathbf{x})} \right)^{-1}$$

Logistic regression

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n -\log P_f(y_i | \mathbf{x}_i)$$

$$\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \left(-\mathbb{I}[y_i = 1] \log \eta(\mathbf{x}_i) - \mathbb{I}[y_i = -1] \log(1 - \eta(\mathbf{x}_i)) \right)$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \left(\mathbb{I}[y_i = 1] \log \left(1 + e^{-f(\mathbf{x}_i)} \right) + \mathbb{I}[y_i = -1] \log \left(1 + e^{f(\mathbf{x}_i)} \right) \right)$$

$$= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right).$$

$$\eta(\mathbf{x}) = \left(1 + e^{-f(\mathbf{x})} \right)^{-1}$$

$$1 - \eta(\mathbf{x}) = 1 - \frac{1}{1 + e^{-f(\mathbf{x})}} = \frac{e^{-f(\mathbf{x})}}{1 + e^{-f(\mathbf{x})}} = \left(1 + e^{f(\mathbf{x})} \right)^{-1}$$

Logistic regression

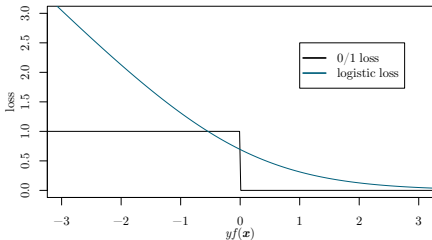
$$\begin{aligned}\hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right) \\ &= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right) .\end{aligned}$$

Logistic regression

$$\begin{aligned}\hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right) \\ &= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right).\end{aligned}$$

- Effectively, we replaced 0/1 loss with **logistic loss**:

$$\ell_{\log}(y, f(\mathbf{x})) = \log \left(1 + e^{-yf(\mathbf{x})} \right).$$

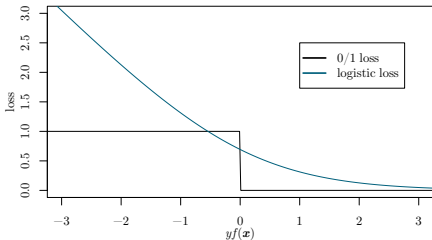


Logistic regression

$$\begin{aligned}\hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right) \\ &= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right).\end{aligned}$$

- Effectively, we replaced 0/1 loss with **logistic loss**:

$$\ell_{\log}(y, f(\mathbf{x})) = \log \left(1 + e^{-yf(\mathbf{x})} \right).$$



- Commonly used, better than least squares in practice.

Learning algorithm for logistic regression

- Let $f(\mathbf{x})$ be a linear function of input attributes:

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^m w_j x_j = \mathbf{w} \cdot \mathbf{x}.$$

- The task of the learning algorithm is to solve:

$$\begin{aligned}\hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \ell_{\log}(y_i, \mathbf{w} \cdot \mathbf{x}_i) \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)).\end{aligned}$$

- This problem is usually solved using iterative convex optimization algorithms.

Learning algorithm for logistic regression

- Consider a simple gradient descent algorithm.
- Total loss over the training examples:

$$\hat{L}(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}))$$

- The gradient descent algorithm:
 - ▶ Initialize \mathbf{w}^0 , e.g. by $\mathbf{w}^0 = \mathbf{0}$
 - ▶ Repeat until convergence:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha \frac{\partial \hat{L}(\mathbf{w}^t)}{\partial \mathbf{w}^t}$$

where α is the step size (learning rate) and

$$\frac{\partial \hat{L}(\mathbf{w})}{\partial \mathbf{w}} = \left(\frac{\partial \hat{L}}{\partial w_1}, \dots, \frac{\partial \hat{L}}{\partial w_m} \right)$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\frac{\partial L}{\partial w_j} =$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\frac{\partial L}{\partial w_j} = - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1}\end{aligned}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1}\end{aligned}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1} \frac{\exp(\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1} \frac{\exp(\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

- Denote:

$$\hat{\eta}_i = \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)}$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1} \frac{\exp(\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

- Denote:

$$\hat{\eta}_i = \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)}$$

- Then:

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1} \frac{\exp(\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

- Denote:

$$\hat{\eta}_i = \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)}$$

- Then:

$$\frac{\partial L}{\partial w_j} = - \sum_{y_i=1} (1 - \hat{\eta}_i) x_{ij} + \sum_{y_i=-1} \hat{\eta}_i x_{ij} =$$

Learning algorithm for logistic regression

- Compute the partial derivative of the loss with respect to w_j :

$$\begin{aligned}\frac{\partial L}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{y_i=1} \frac{\exp(-\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)} + \sum_{y_i=-1} \frac{\exp(\mathbf{w} \cdot \mathbf{x}_i) x_{ij}}{1 + \exp(\mathbf{w} \cdot \mathbf{x}_i)}\end{aligned}$$

- Denote:

$$\hat{\eta}_i = \frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}_i)}$$

- Then:

$$\frac{\partial L}{\partial w_j} = - \sum_{y_i=1} (1 - \hat{\eta}_i) x_{ij} + \sum_{y_i=-1} \hat{\eta}_i x_{ij} = - \sum_{i=1}^n (y'_i - \hat{\eta}_i) x_{ij},$$

where $y' = \frac{y+1}{2} \in \{0, 1\}$.

Learning algorithm for logistic regression

- Connection with linear regression:
 - ▶ The partial derivative of the squared loss over training examples with respect to w_j is similar:

$$\frac{\partial \hat{L}_{se}}{\partial w_j} =$$

Learning algorithm for logistic regression

- Connection with linear regression:
 - ▶ The partial derivative of the squared loss over training examples with respect to w_j is similar:

$$\frac{\partial \hat{L}_{se}}{\partial w_j} = -2 \sum_{i=1}^n (y_i - f(\mathbf{x})) x_{ij}$$

Learning algorithm for logistic regression

- Connection with linear regression:
 - ▶ The partial derivative of the squared loss over training examples with respect to w_j is similar:

$$\frac{\partial \hat{L}_{se}}{\partial w_j} = -2 \sum_{i=1}^n (y_i - f(\mathbf{x})) x_{ij}$$

- ▶ For the squared error loss, however, the solution can be found analytically since $f(\mathbf{x})$ is linear.

Learning algorithm for logistic regression

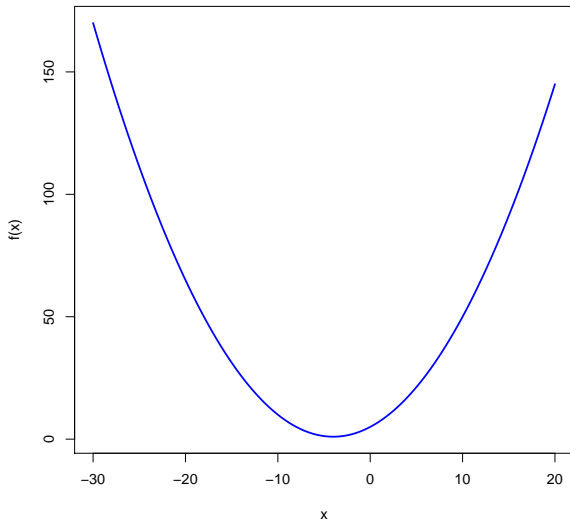
- Connection with linear regression:
 - ▶ The partial derivative of the squared loss over training examples with respect to w_j is similar:

$$\frac{\partial \hat{L}_{se}}{\partial w_j} = -2 \sum_{i=1}^n (y_i - f(\mathbf{x})) x_{ij}$$

- ▶ For the squared error loss, however, the solution can be found analytically since $f(\mathbf{x})$ is linear.
 - ▶ For logistic regression we need to use iterative methods, since $\hat{\eta}_i$ is not linear, but sigmoid.
- The simplest method assumes α to be constant.
- It does not mean that the step is always the same:
 - ▶ As we approach (local) minimum, $\frac{\partial \hat{L}(\mathbf{w})}{\partial \mathbf{w}}$ takes smaller values.

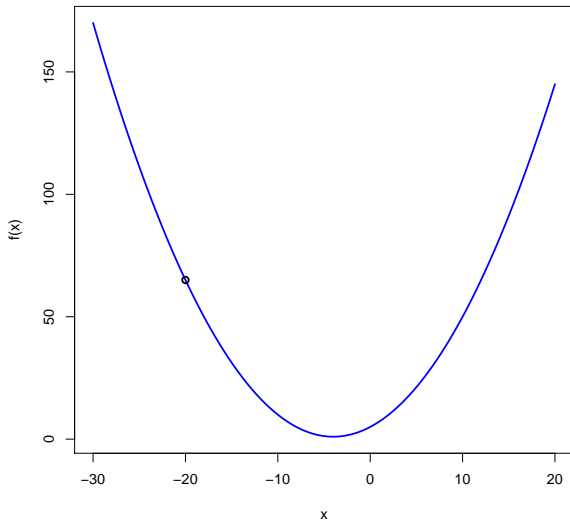
Logistic regression

- Gradient descent example:



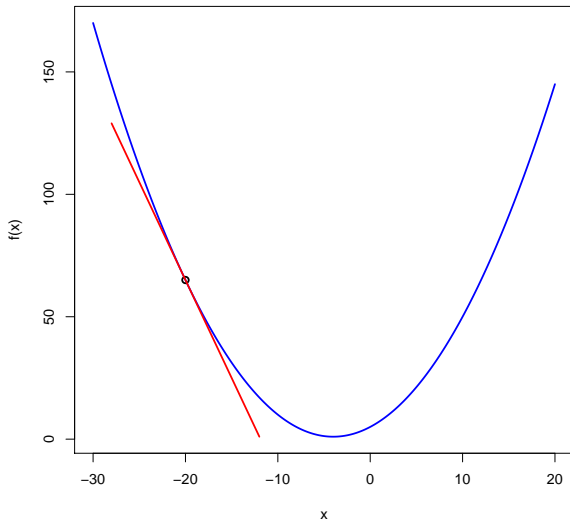
Logistic regression

- Gradient descent example:



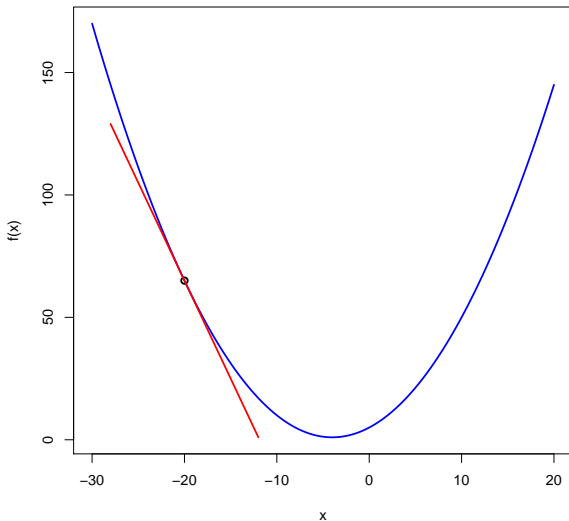
Logistic regression

- Gradient descent example:



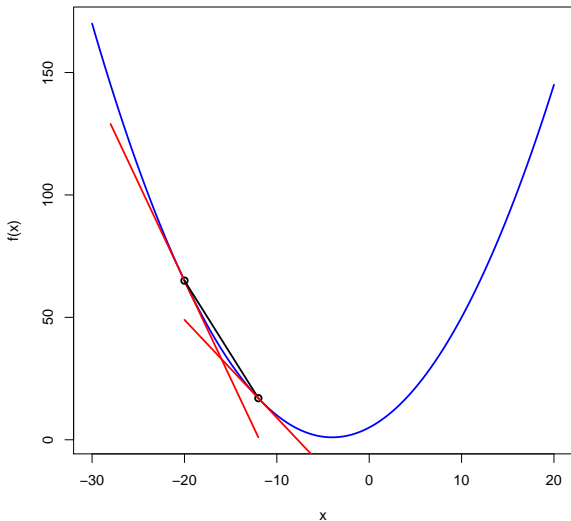
Logistic regression

- Gradient descent example ($\alpha = 1$):



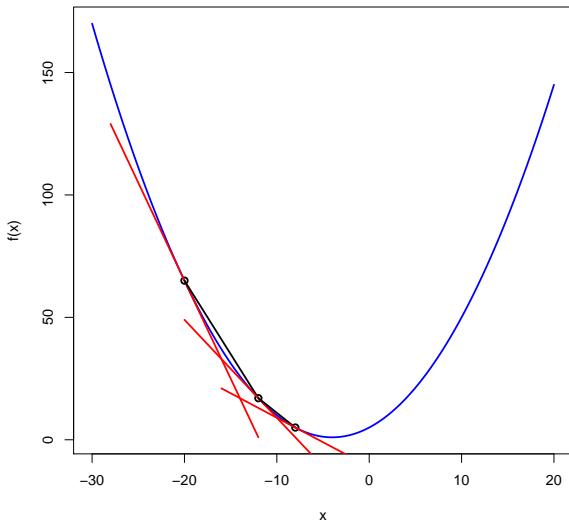
Logistic regression

- Gradient descent example ($\alpha = 1$):



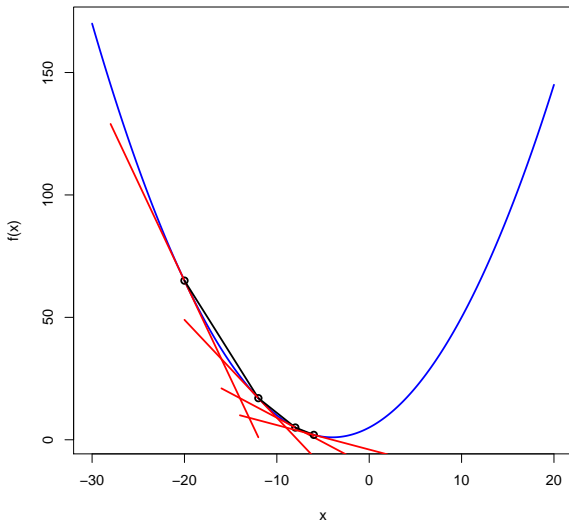
Logistic regression

- Gradient descent example ($\alpha = 1$):



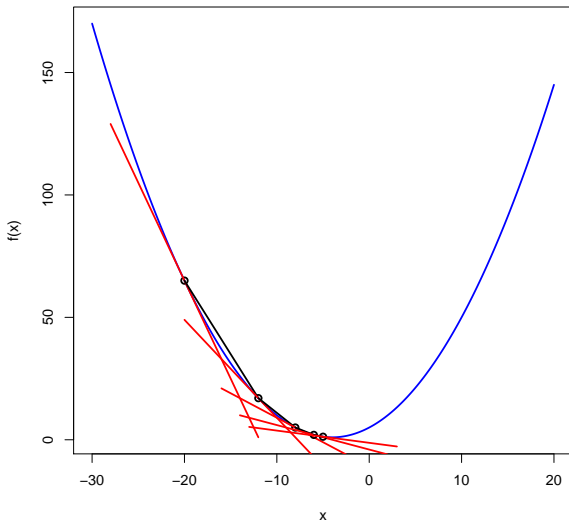
Logistic regression

- Gradient descent example ($\alpha = 1$):



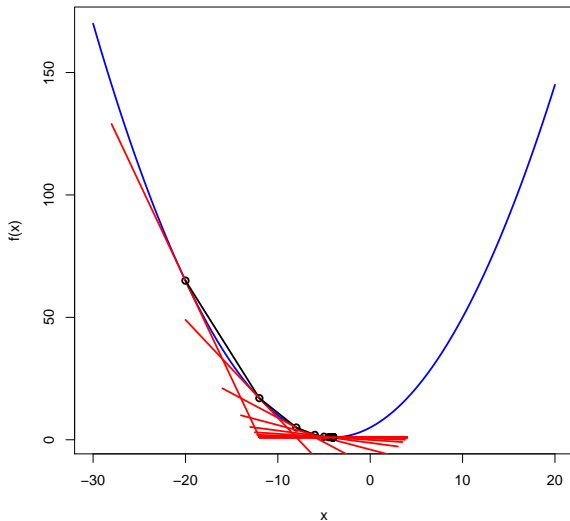
Logistic regression

- Gradient descent example ($\alpha = 1$):



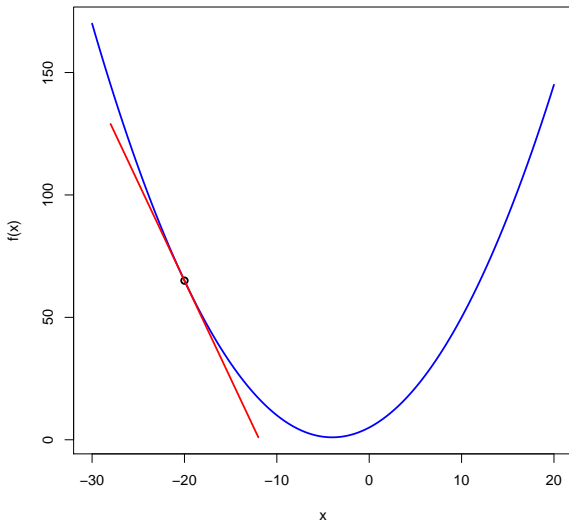
Logistic regression

- Gradient descent example ($\alpha = 1$):



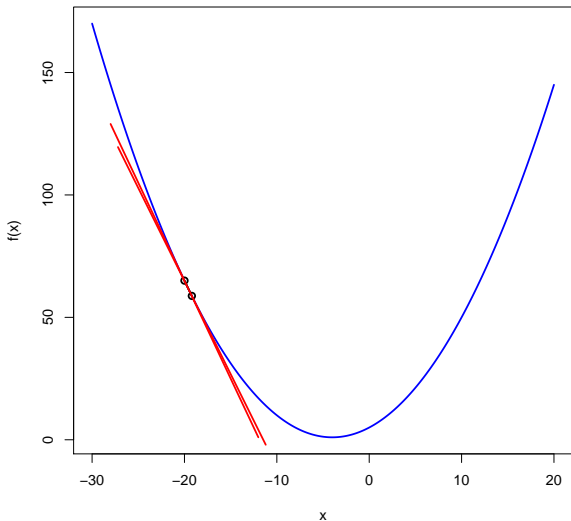
Logistic regression

- Gradient descent example($\alpha = 0.1$):



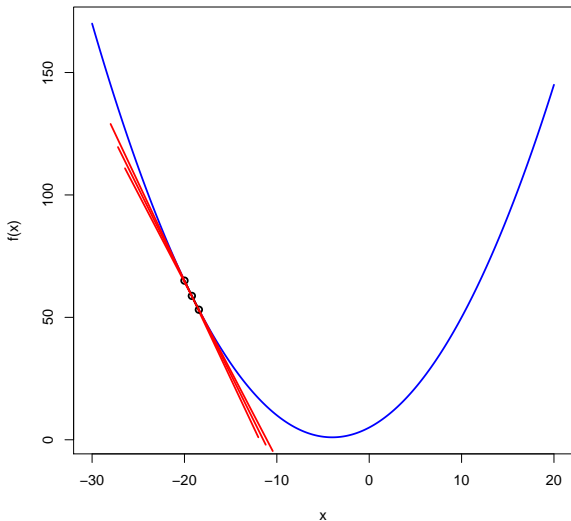
Logistic regression

- Gradient descent example($\alpha = 0.1$):



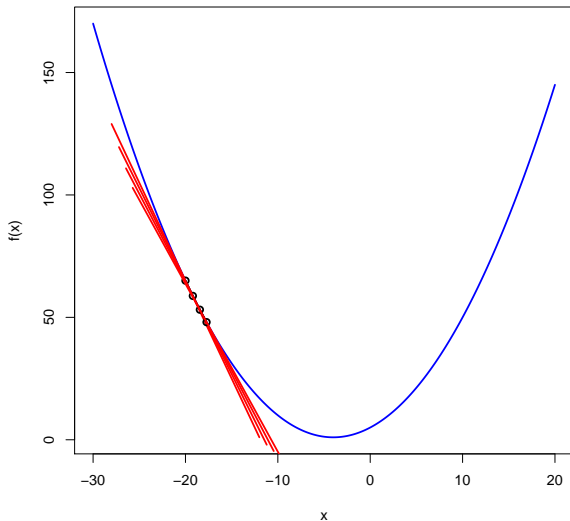
Logistic regression

- Gradient descent example($\alpha = 0.1$):



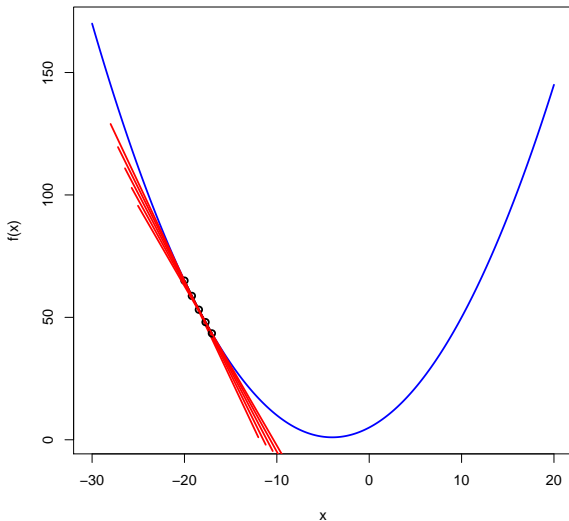
Logistic regression

- Gradient descent example($\alpha = 0.1$):



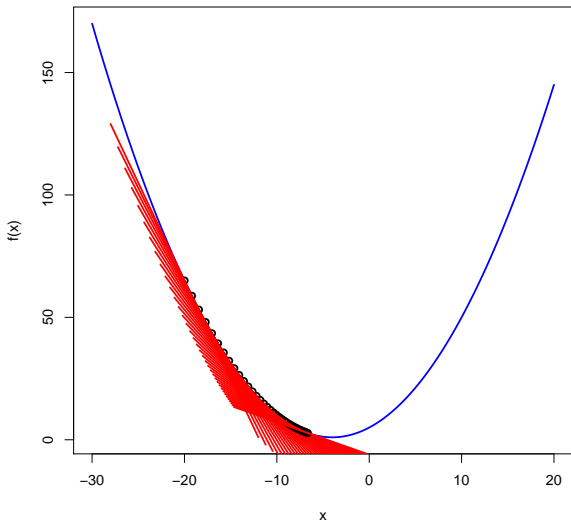
Logistic regression

- Gradient descent example($\alpha = 0.1$):



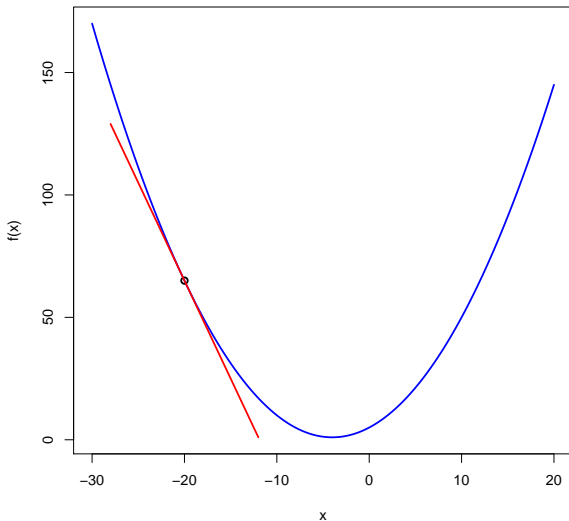
Logistic regression

- Gradient descent example($\alpha = 0.1$):



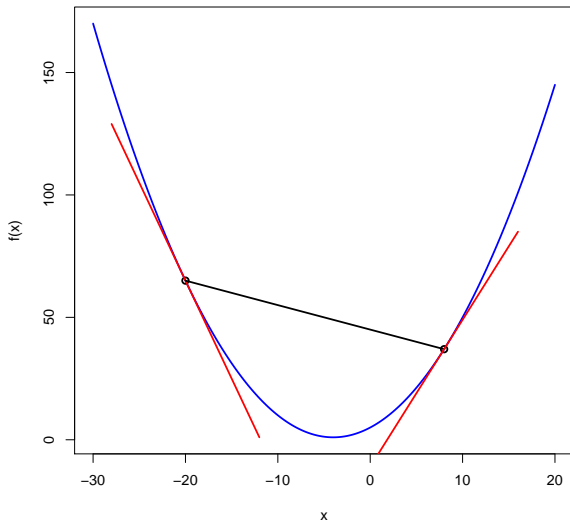
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



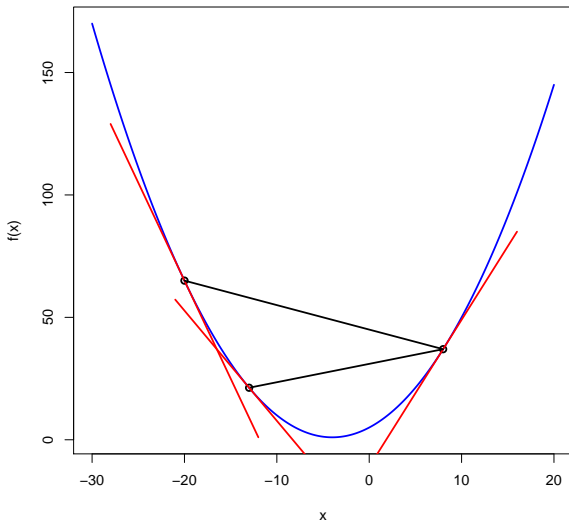
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



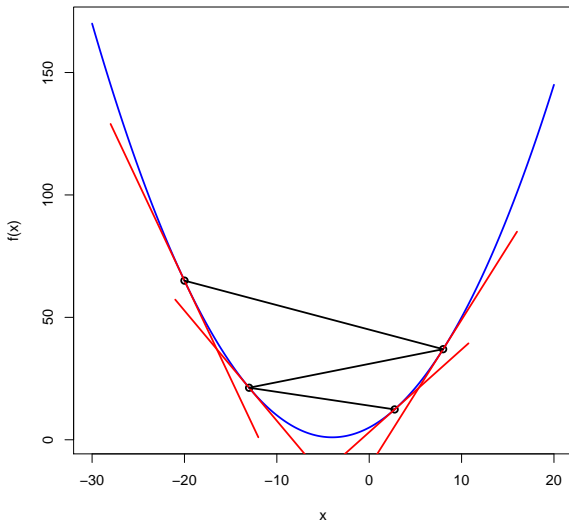
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



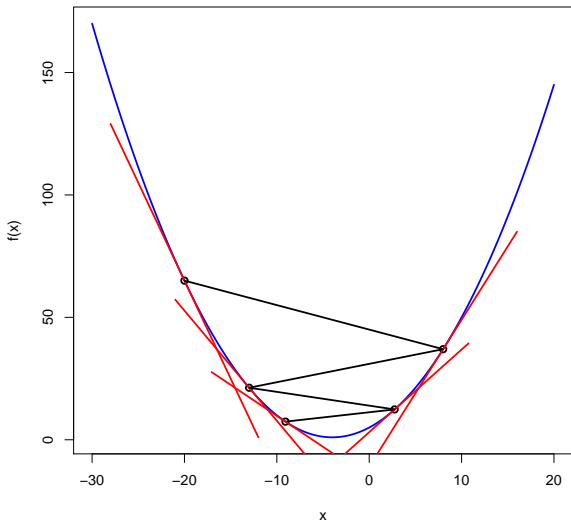
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



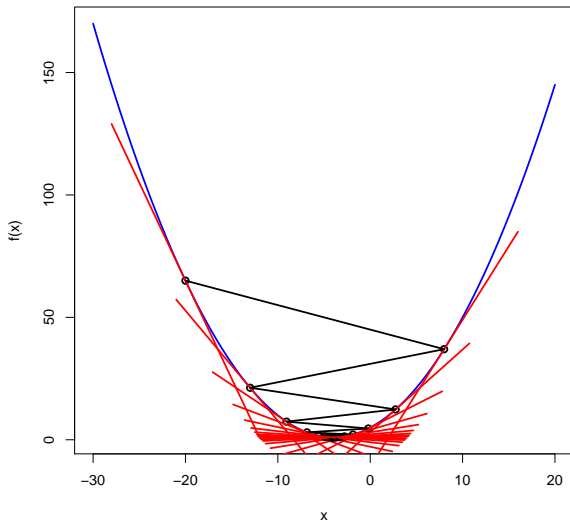
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



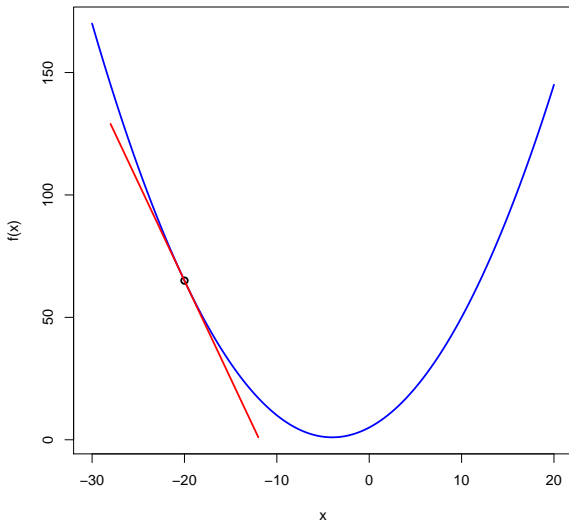
Logistic regression

- Gradient descent example ($\alpha = 0.1$):



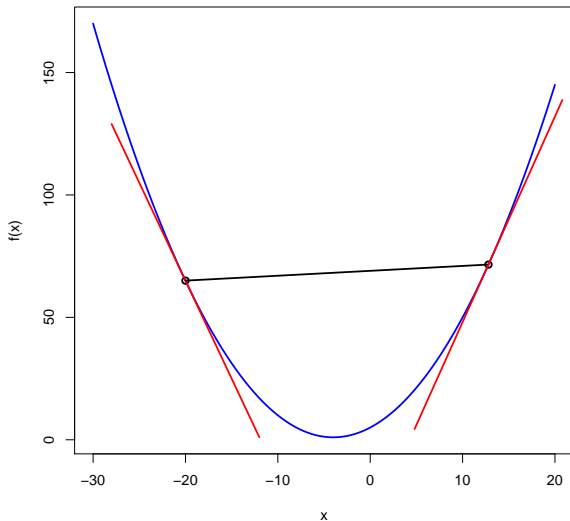
Logistic regression

- Gradient descent example ($\alpha = 4.1$):



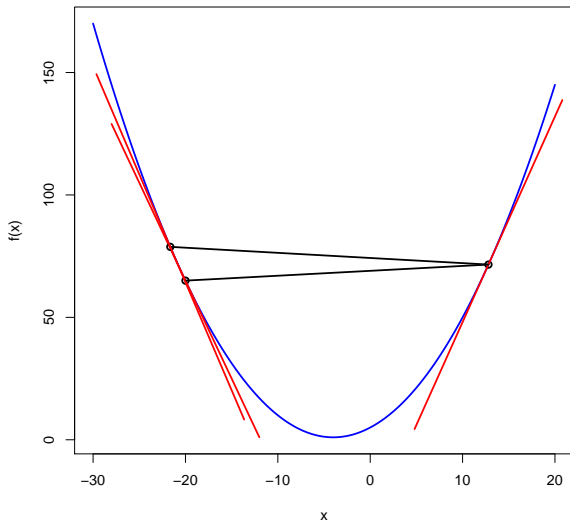
Logistic regression

- Gradient descent example ($\alpha = 4.1$):



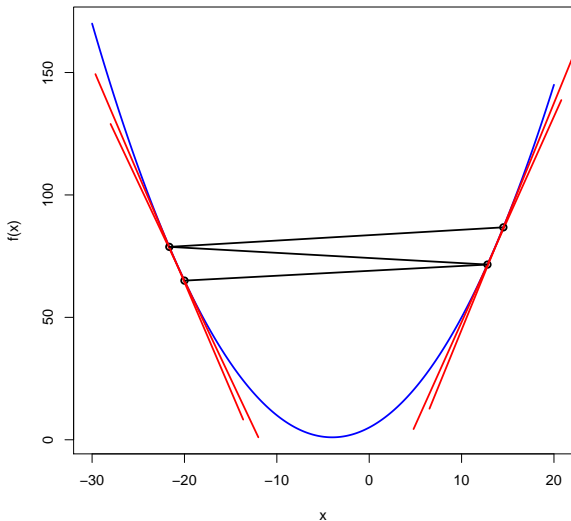
Logistic regression

- Gradient descent example ($\alpha = 4.1$):



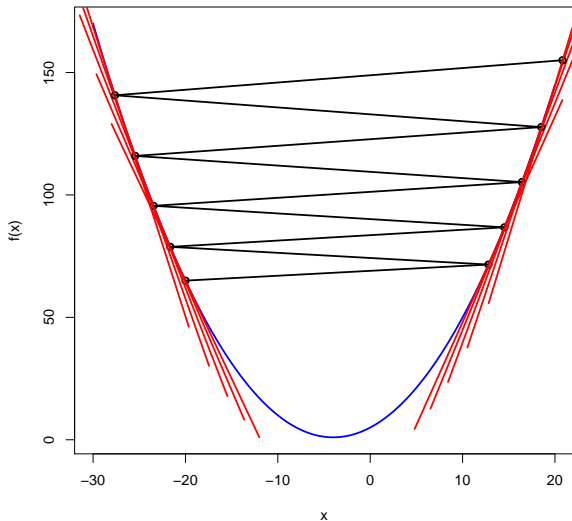
Logistic regression

- Gradient descent example ($\alpha = 4.1$):



Logistic regression

- Gradient descent example ($\alpha = 4.1$):



Stochastic gradient descent

- Computing the gradient can be costly:
 - ▶ Sum over n components.
 - ▶ The number of training examples n can be large.
- We can approximate the gradient by sampling from the training set.
- **Stochastic gradient descent:**
 - ▶ Randomly draw an example from the training set.
 - ▶ Compute the gradient based on this single example:

$$\frac{\partial \hat{L}_i}{\partial w_j} = -(y'_i - \hat{\eta}_i)x_{ij} ,$$

where $\hat{L}_i = \ell_{\log}(y_i, \mathbf{w} \cdot \mathbf{x}_i)$.

- ▶ Update parameters.
- ▶ Repeat until convergence.

Stochastic gradient descent

```
Input: learning rate  $\alpha$ 
 $w = \mathbf{0}$ ; //(or use random values)
while (approximate minimum is obtained) {
    Randomly shuffle examples in the training set
    for  $i=1$  to  $n$  {
         $w := w - \alpha \frac{\partial \hat{L}_i(w)}{\partial w}$ 
    }
}
```

Loss functions

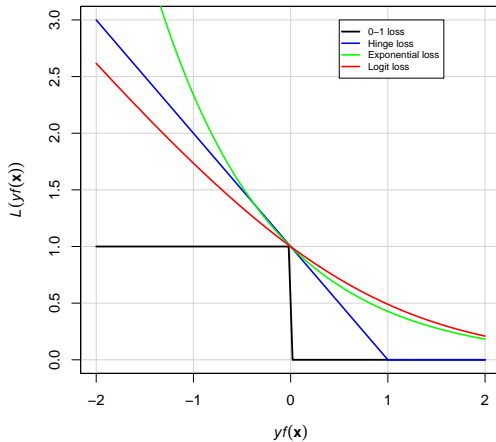


Figure: Loss functions for classification task

Outline

- 1 Linear Models for Classification
- 2 Summary

Summary

- Linear models for classification:
 - ▶ Linear regression.
 - ▶ Perceptron.
 - ▶ Support vector machines.
 - ▶ Logistic regression.

Bibliography

- T. Hastie, R. Tibshirani, and J. Friedman. *Elements of Statistical Learning: Second Edition*.
Springer, 2009
<http://www-stat.stanford.edu/~tibs/ElemStatLearn/>
- Christopher M. Bishop. *Pattern Recognition and Machine Learning*.
Springer-Verlag, 2006
- David Barber. *Bayesian Reasoning and Machine Learning*.
Cambridge University Press, 2012
<http://www.cs.ucl.ac.uk/staff/d.barber/brml/>
- Yaser S. Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin. *Learning From Data*.
AMLBook, 2012