

Decision-theoretic Machine Learning

Krzysztof Dembczyński and Wojciech Kotłowski

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



Poznań University of Technology, Summer 2019

Agenda

- 1 Introduction to Machine Learning
- 2 **Binary Classification**
- 3 Bipartite Ranking
- 4 Multi-Label Classification

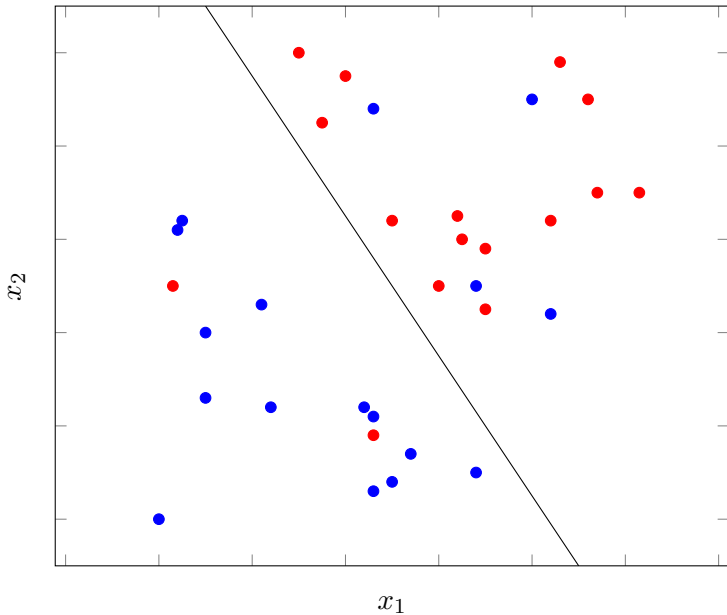
Outline

- 1 Linear models for classification
- 2 Some statistical decision theory
- 3 Classification calibrated losses
- 4 Generalization error

Outline

- 1 Linear models for classification
- 2 Some statistical decision theory
- 3 Classification calibrated losses
- 4 Generalization error

Linear models for classification



Linear models for classification

- Let the output variable be $y \in \{-1, 1\}$ (or alternatively $y \in \{0, 1\}$).
- Prediction function $h(\mathbf{x}) \in \{-1, 1\}$ (or alternatively $h(\mathbf{x}) \in \{0, 1\}$).

Linear models for classification

- Let the output variable be $y \in \{-1, 1\}$ (or alternatively $y \in \{0, 1\}$).
- Prediction function $h(\mathbf{x}) \in \{-1, 1\}$ (or alternatively $h(\mathbf{x}) \in \{0, 1\}$).
- 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

- Let the output variable be $y \in \{-1, 1\}$ (or alternatively $y \in \{0, 1\}$).
- Prediction function $h(\mathbf{x}) \in \{-1, 1\}$ (or alternatively $h(\mathbf{x}) \in \{0, 1\}$).
- 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

- Let the output variable be $y \in \{-1, 1\}$ (or alternatively $y \in \{0, 1\}$).
- Prediction function $h(\mathbf{x}) \in \{-1, 1\}$ (or alternatively $h(\mathbf{x}) \in \{0, 1\}$).
- 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** $f \in \mathcal{F}$:

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

- ▶ **Threshold** f at 0:

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

so that $\llbracket y_i h(\mathbf{x}_i) \leq 0 \rrbracket = \llbracket y_i f(\mathbf{x}_i) \leq 0 \rrbracket$.

Linear models for classification

- Solve:

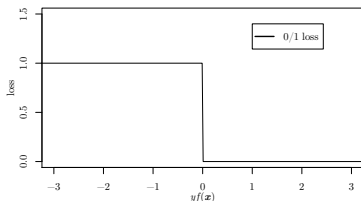
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{I}[y_i f(\mathbf{x}_i) \leq 0]$$

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**.

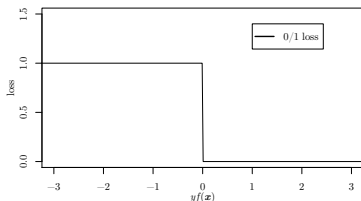


Linear models for classification

- Solve:

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

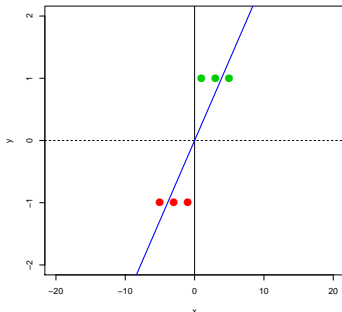
- Still **hard** to optimize: 0/1 loss is **discontinuous** and **non-convex**.
 - ▶ e.g., when \mathcal{H} is a class of linear function, the problem known to be **NP-hard**.



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

Linear regression for classification

- Let us assume that $y \in \{-1, 1\}$.
- We can try to use linear regression to solve binary classification problem:

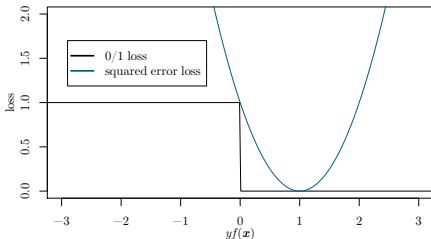


- Use a threshold to decide which class is predicted:
Positive class: $f(\mathbf{x}) > 0$ Negative class: $f(\mathbf{x}) < 0$
- Minimization of squared loss is not bad, since it leads to estimation of the conditional probability, i.e., $P(y = 1|\mathbf{x}) = \frac{1}{2}\mathbb{E}(y|\mathbf{x}) + \frac{1}{2}$.

Linear regression for classification

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

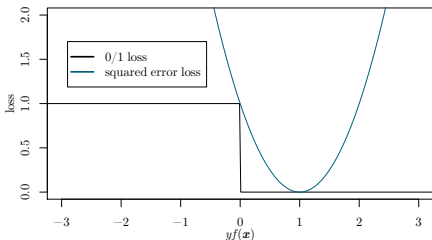
$$\ell_{\text{sq}}(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2 = (1 - yf(\mathbf{x}))^2.$$



Linear regression for classification

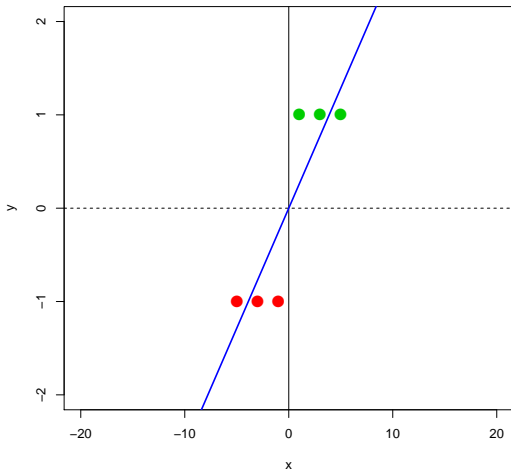
- Effectively, we replaced 0/1 loss by **squared error loss**:

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

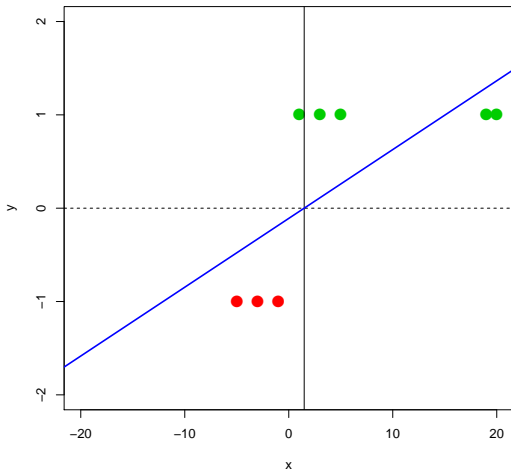


- Works nicely in practice, but has several drawbacks ...

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

- One possibility is to consider only those examples that are incorrectly classified.
- The quantity $yf(\mathbf{x})$ is usually referred to as **margin**.
- If y and $f(\mathbf{x})$ have the same sign (i.e., the prediction is correct), then

$$yf(\mathbf{x}) > 0,$$

otherwise

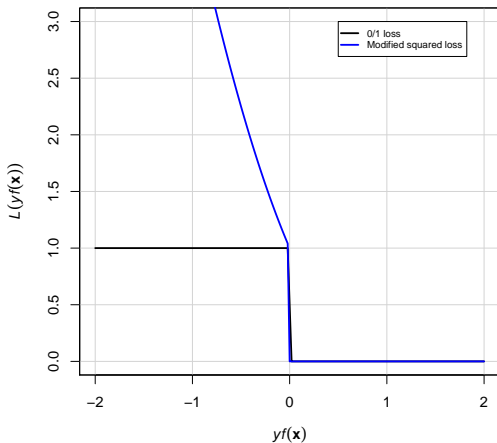
$$yf(\mathbf{x}) \leq 0.$$

Linear models for classification

- We could consider the following loss function:

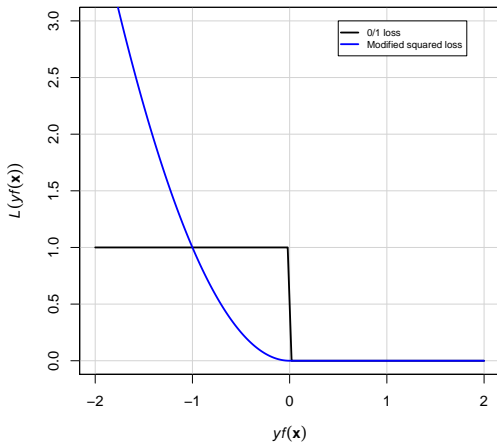
$$\ell(y, f(\mathbf{x})) = \begin{cases} 0, & \text{if } yf(\mathbf{x}) > 0 \\ (1 - yf(\mathbf{x}))^2, & \text{otherwise.} \end{cases}$$

Linear models for classification



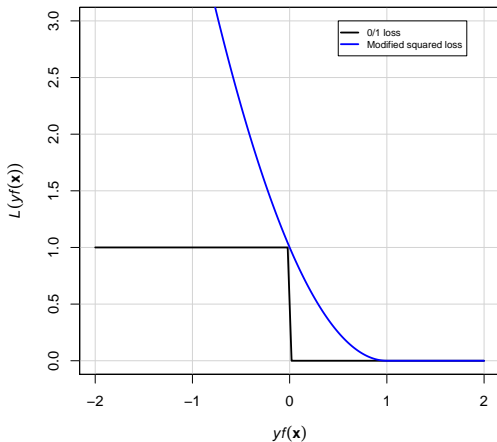
- This definition does not lead to a “nice” shape of the loss.

Linear models for classification



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

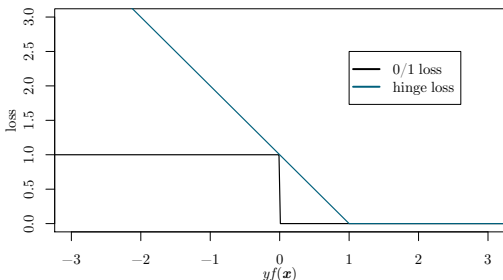
Linear models for classification



- 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:

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

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 the stochastic gradient descent algorithm.

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  $n$  {
         $w := w - \alpha \frac{\partial \ell(y_i, f(\mathbf{x}_i))}{\partial w}$ 
    }
}
```

Perceptron

- Learning algorithm for perceptron:
 - ▶ 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

- Learning algorithm for perceptron:
 - ▶ 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:

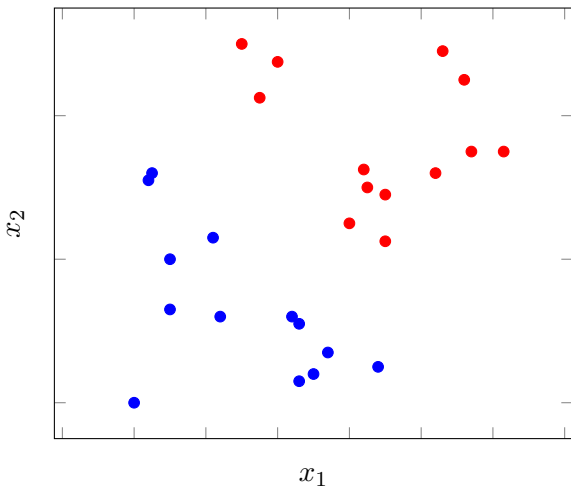
$$\frac{\partial \ell_{\text{perc}}(y, f(\mathbf{x}))}{\partial \mathbf{w}} = -y\mathbf{x} .$$

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

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

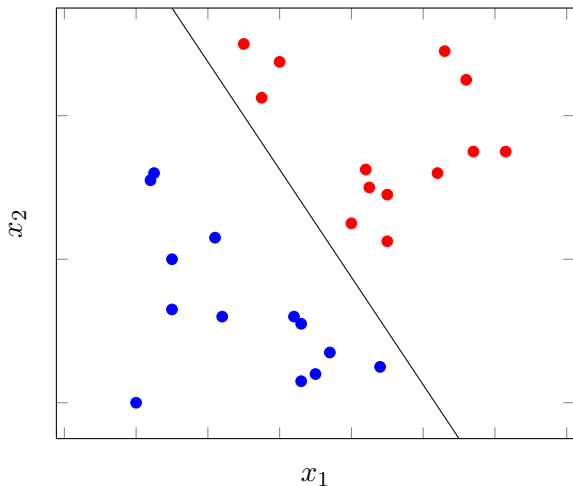
Perceptron – Graphical interpretation

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



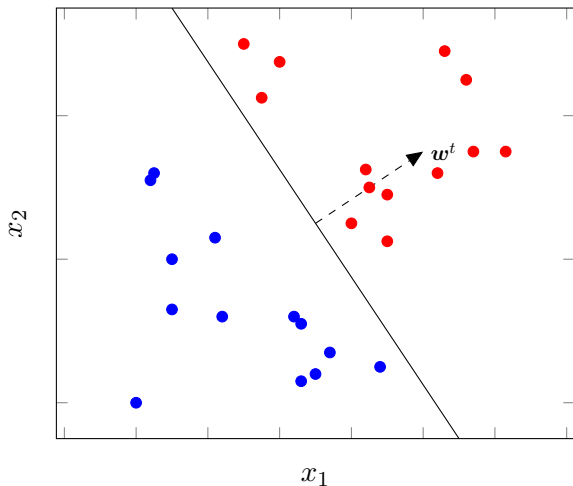
Perceptron – Graphical interpretation

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



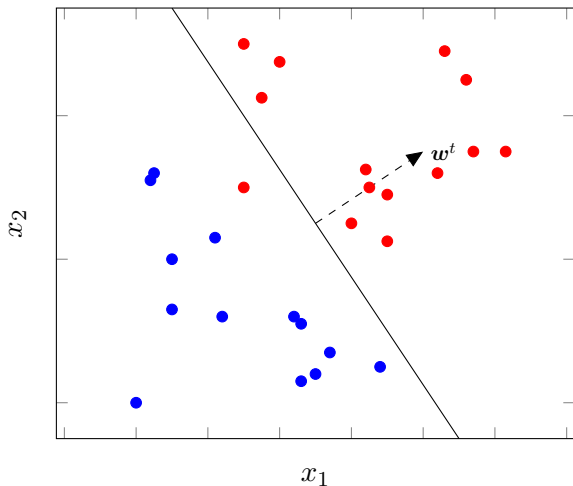
Perceptron – Graphical interpretation

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



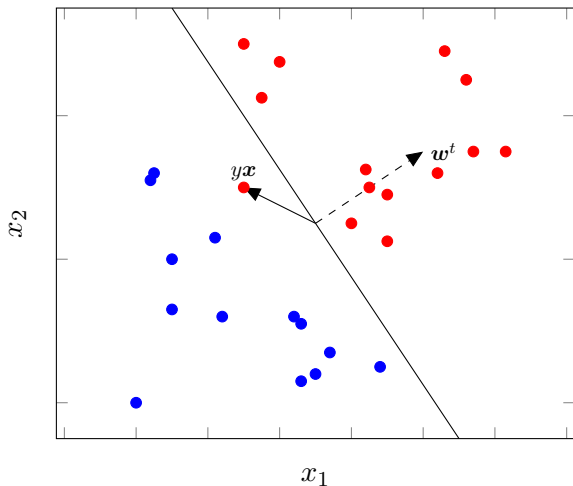
Perceptron – Graphical interpretation

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



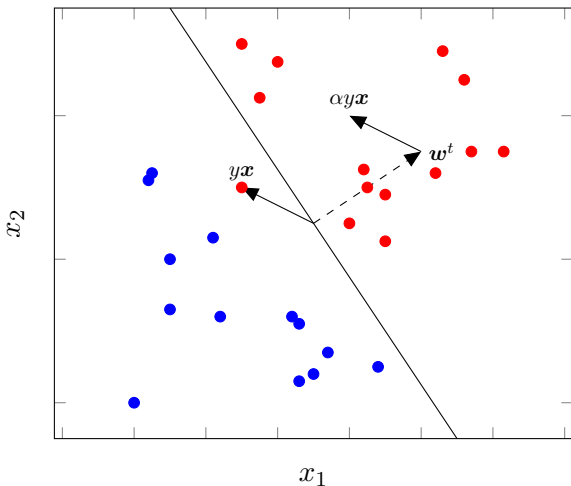
Perceptron – Graphical interpretation

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



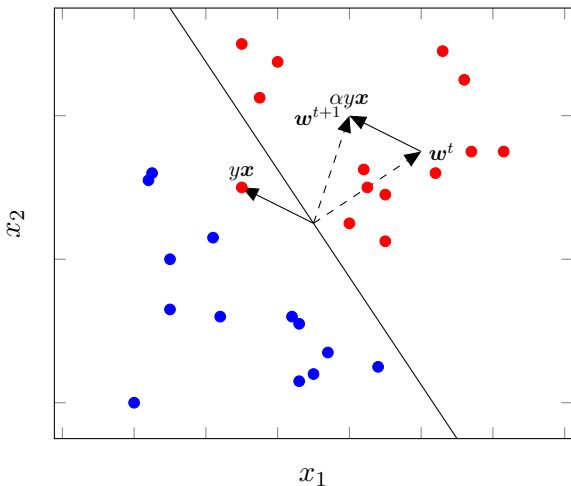
Perceptron – Graphical interpretation

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



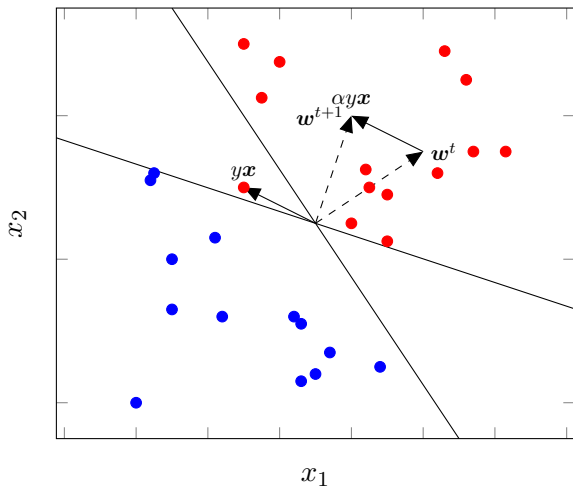
Perceptron – Graphical interpretation

- 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:



Linear models with hinge loss

- Loss function:

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

Linear models with hinge loss

- Loss function:

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

- Linear model:

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

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:

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

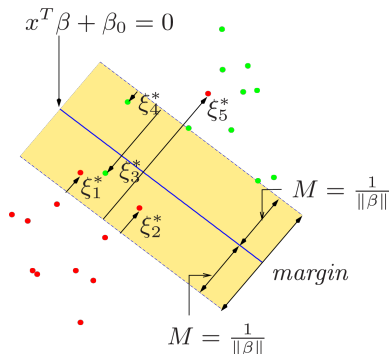
Binary classification by Support Vector Machines (SVM)¹

Find **maximal margin classifier**

$$\min_{\mathbf{w}} \quad \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad y_i \mathbf{w} \cdot \mathbf{x} > 1 - \xi_i \quad \forall i = 1..n$$

$$\xi_i \geq 0 \quad \forall i = 1..n$$



¹ Bernhard E. Boser, Isabelle Guyon, and Vladimir Vapnik. A training algorithm for optimal margin classifiers. In *Conference on Learning Theory (COLT)*, pages 144–152, 1992

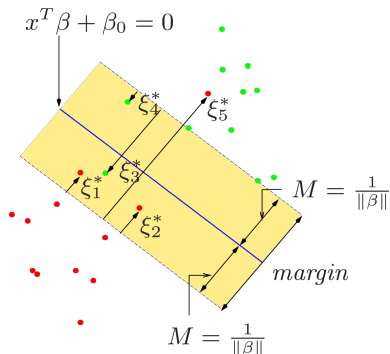
Binary classification by Support Vector Machines (SVM)¹

Find **maximal margin classifier**

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

\Leftrightarrow

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

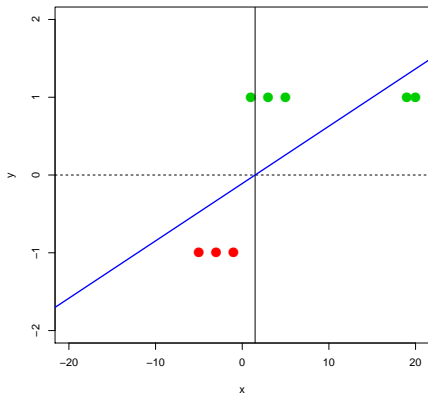


¹ Bernhard E. Boser, Isabelle Guyon, and Vladimir Vapnik. A training algorithm for optimal margin classifiers. In *Conference on Learning Theory (COLT)*, pages 144–152, 1992

Logistic regression

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

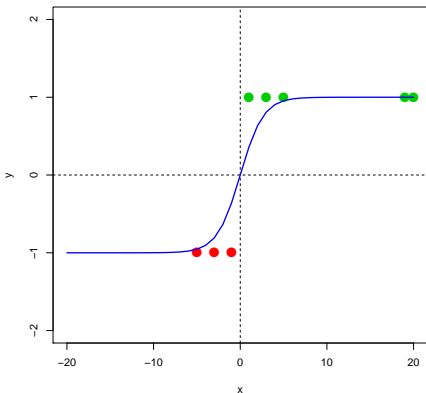
$$g(\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))} \in (0, 1) \quad g(\mathbf{x}) = \frac{1 - \exp(-f(\mathbf{x}))}{1 + \exp(-f(\mathbf{x}))} \in (-1, 1)$$



Logistic regression

- 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 $\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$ 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 $f(\mathbf{x}) \in \mathbb{R}$.

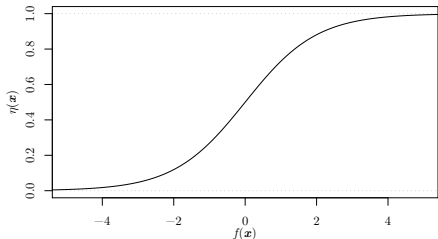
Logistic regression – motivation

- Get an estimate of $\eta(\mathbf{x}) = P(y = 1 | \mathbf{x})$ 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

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



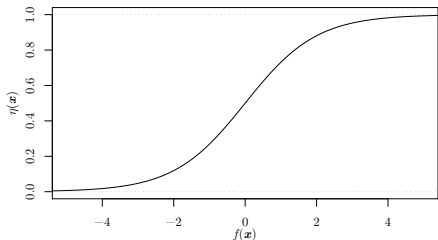
Logistic regression – motivation

- 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**

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

↕

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



Question

Show that if examples are generated by $y \sim P(y)$ and $\mathbf{x}|y \sim N(\mu_y|\Sigma)$ (with shared covariance matrix), then $\log \frac{\eta(\mathbf{x})}{1-\eta(\mathbf{x})}$ is a linear function of \mathbf{x} .

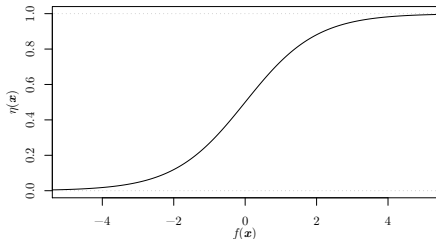
Logistic regression – motivation

- 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**

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

\Updownarrow

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



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

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

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

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

Logistic regression

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

Logistic regression

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

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

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

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

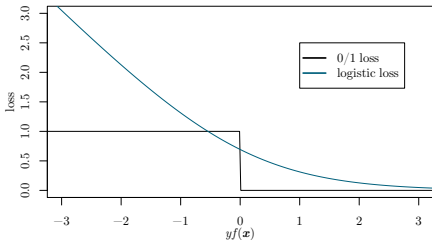
$$\begin{aligned}\hat{f} &= \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right) \\ &= \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i f(\mathbf{x}_i)} \right).\end{aligned}$$

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**:

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

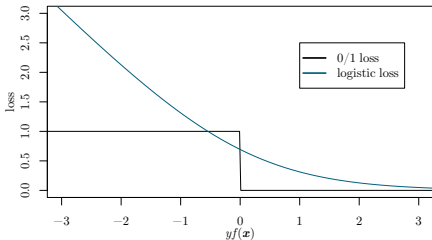


Logistic regression

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

- 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 the 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.

Logistic regression

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

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

- 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

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

Logistic regression

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

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

Logistic regression

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

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

Logistic regression

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

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

Logistic regression

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

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

Logistic regression

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

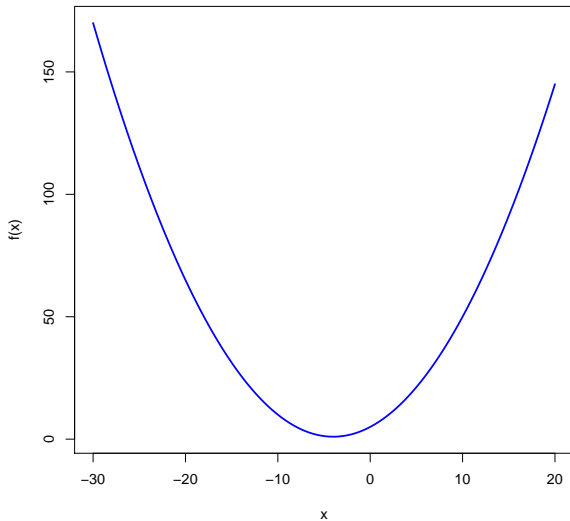
$$\begin{aligned}\frac{\partial \hat{L}}{\partial w_j} &= - \sum_{i=1}^n \frac{\exp(-y_i \mathbf{w} \cdot \mathbf{x}_i) y_i x_{ij}}{1 + \exp(-y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{i=1}^n \frac{y_i x_{ij}}{1 + \exp(y_i \mathbf{w} \cdot \mathbf{x}_i)} \\ &= - \sum_{i=1}^n \beta_i y_i x_{ij},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

- Gradient:

$$\frac{\partial \hat{L}(\mathbf{w})}{\partial \mathbf{w}} = - \sum_{i=1}^n \beta_i y_i \mathbf{x}_i, \quad \text{where } \beta_i = \frac{1}{1 + \exp(y_i \mathbf{w} \cdot \mathbf{x}_i)}$$

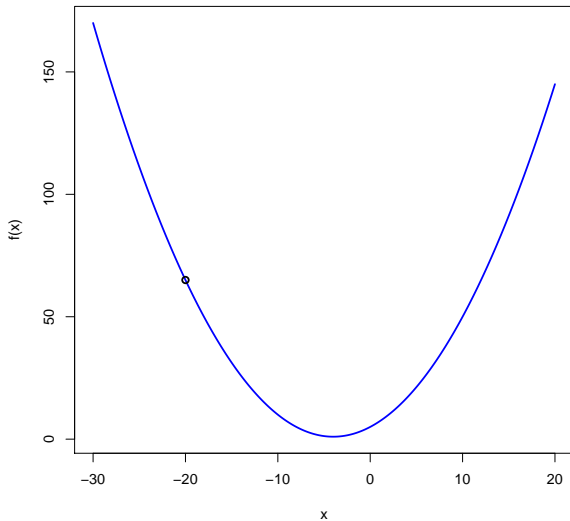
Logistic regression

- Gradient descent example:



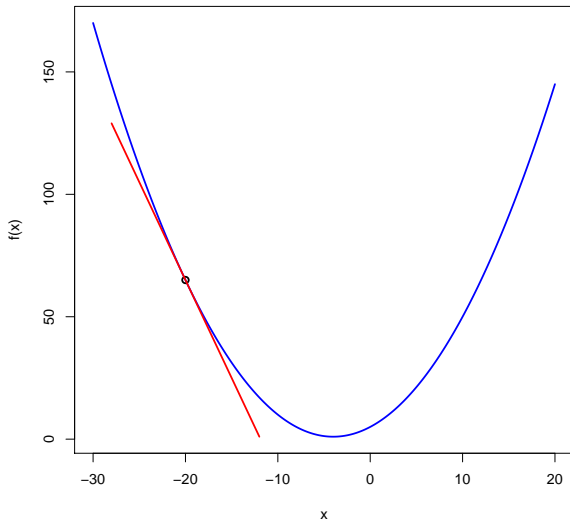
Logistic regression

- Gradient descent example:



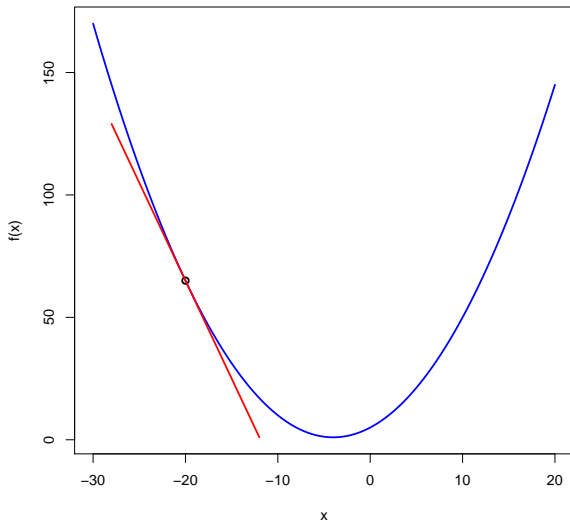
Logistic regression

- Gradient descent example:



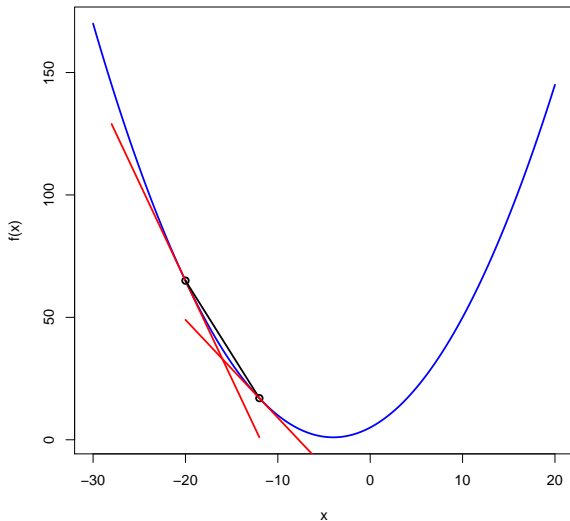
Logistic regression

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



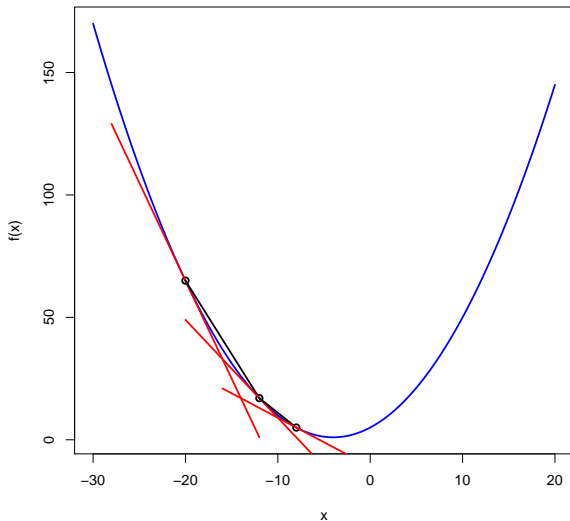
Logistic regression

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



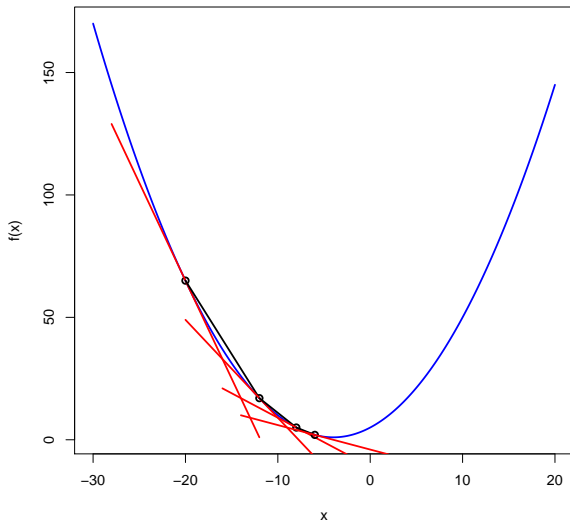
Logistic regression

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



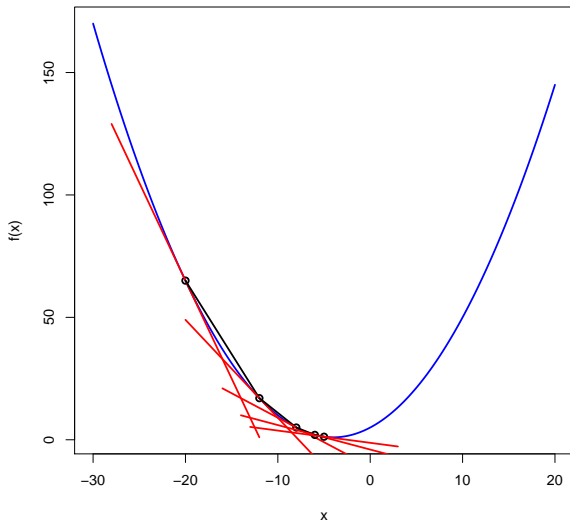
Logistic regression

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



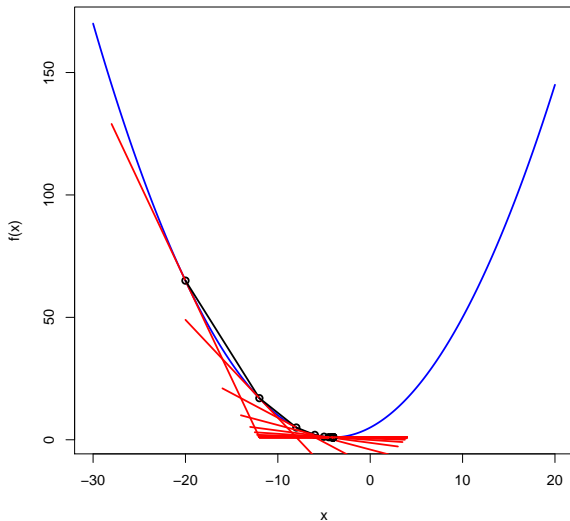
Logistic regression

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



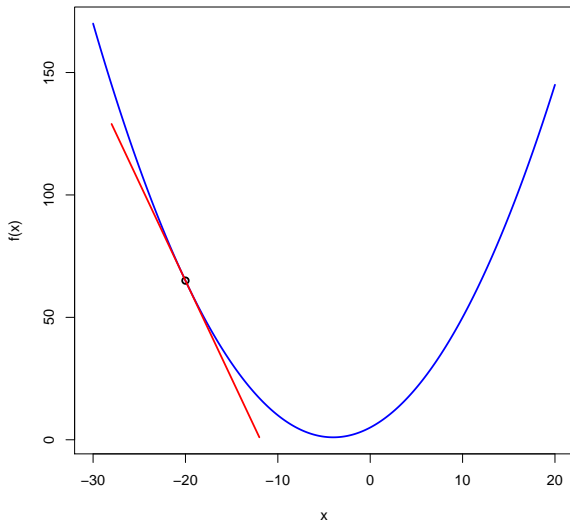
Logistic regression

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



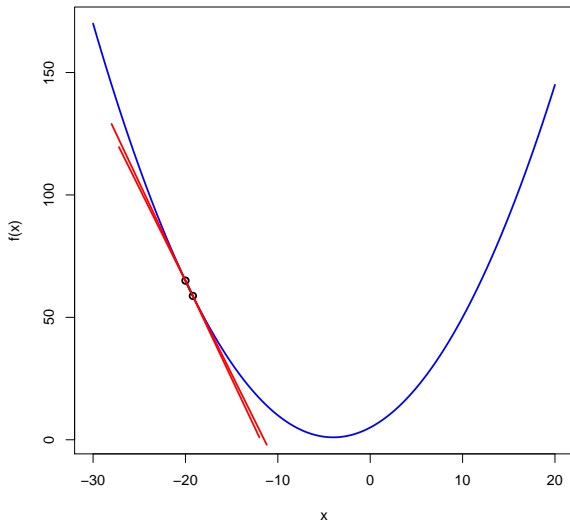
Logistic regression

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



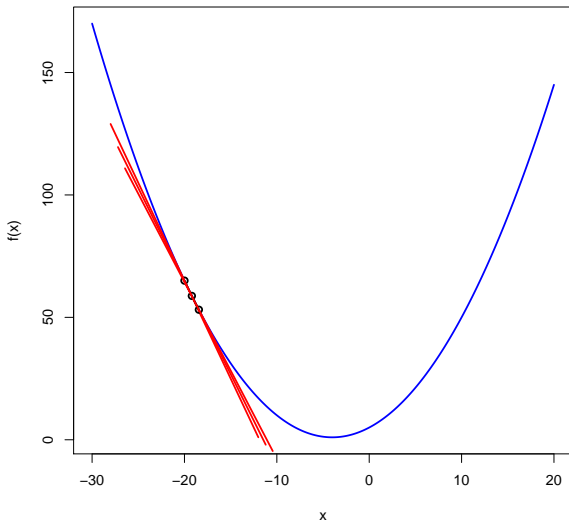
Logistic regression

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



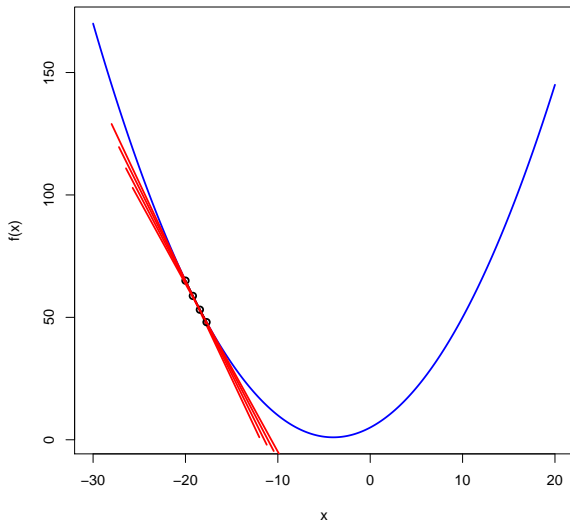
Logistic regression

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



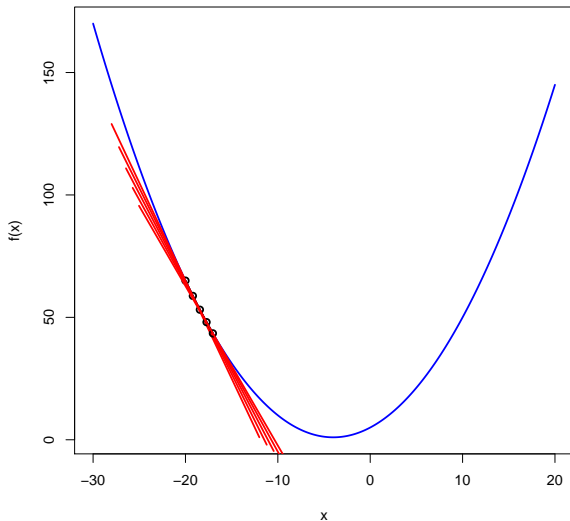
Logistic regression

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



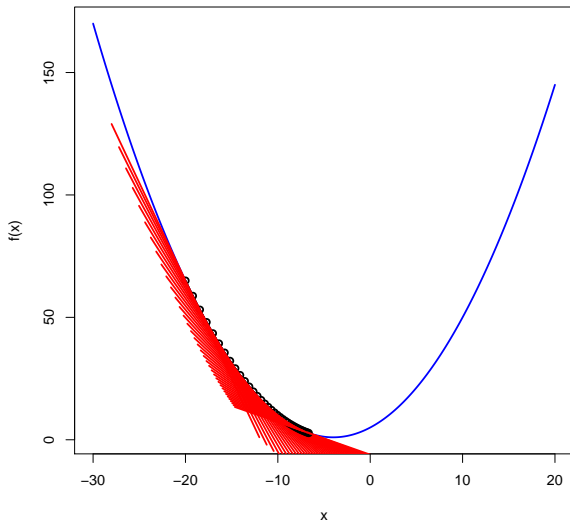
Logistic regression

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



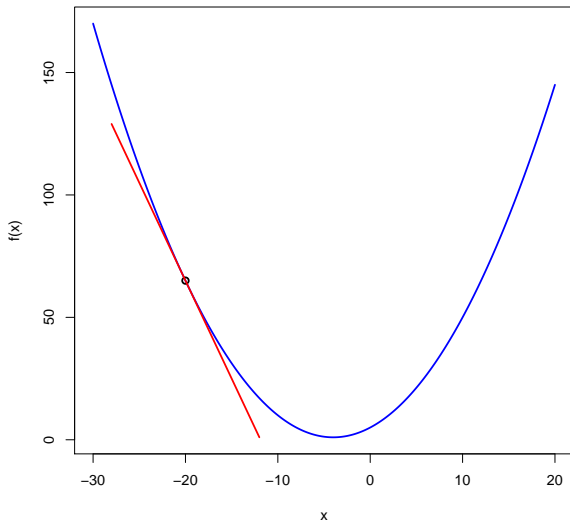
Logistic regression

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



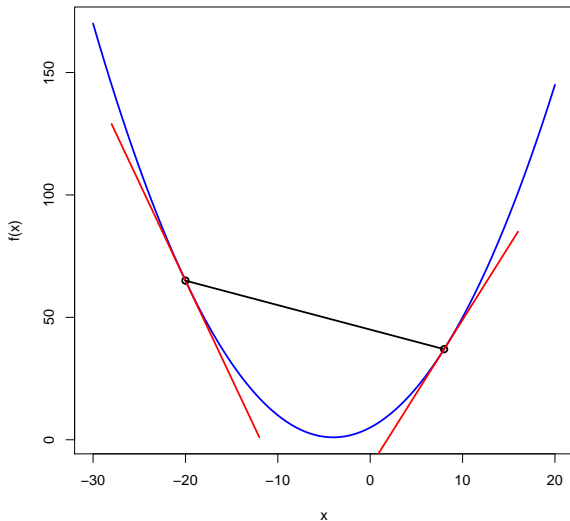
Logistic regression

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



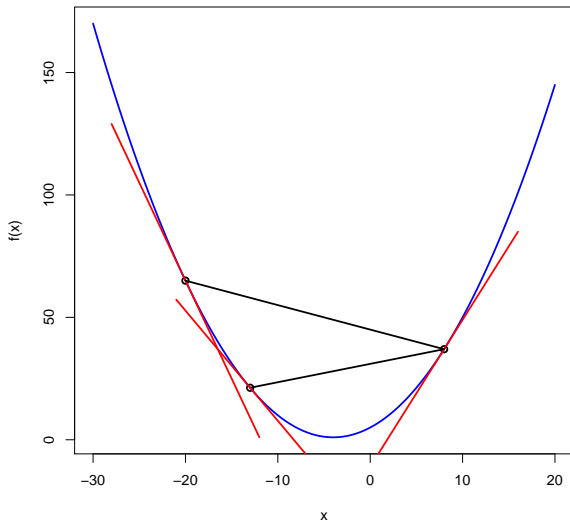
Logistic regression

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



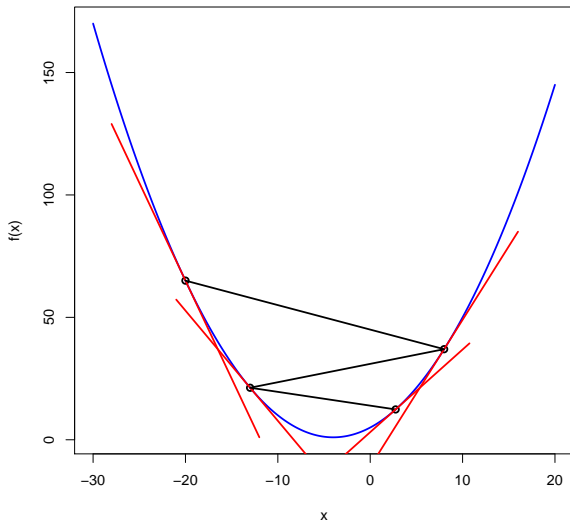
Logistic regression

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



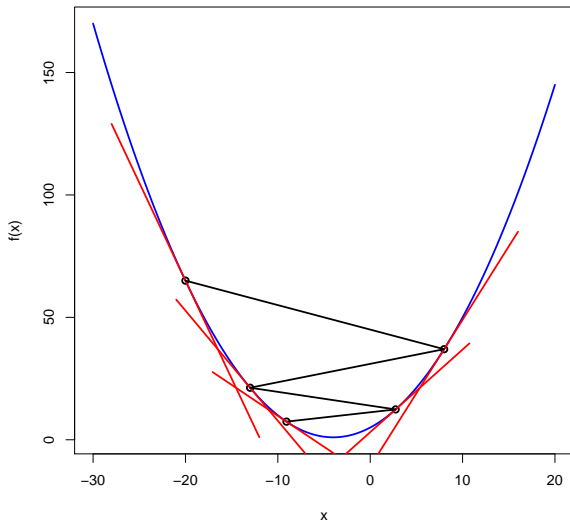
Logistic regression

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



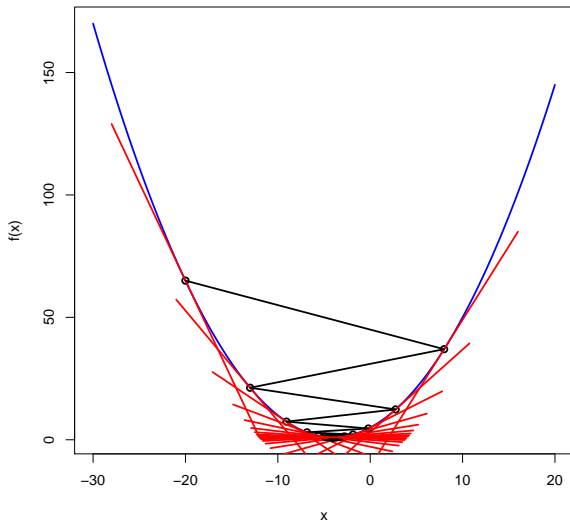
Logistic regression

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



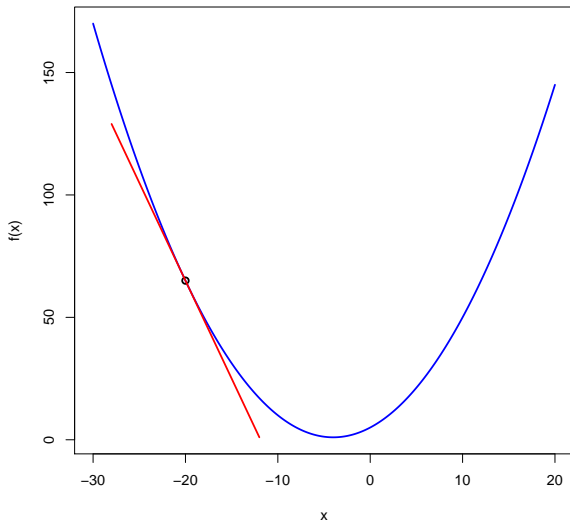
Logistic regression

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



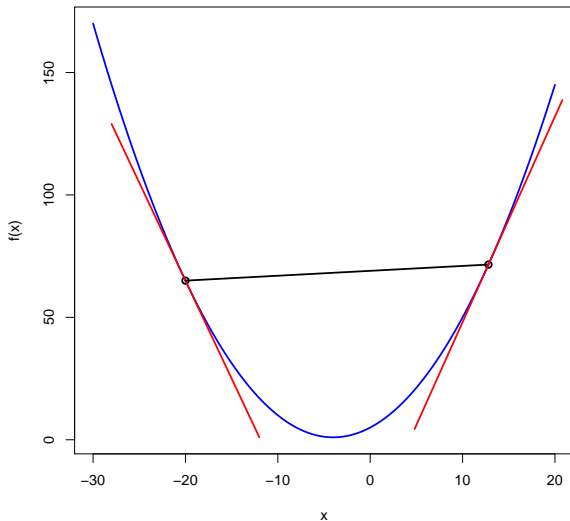
Logistic regression

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



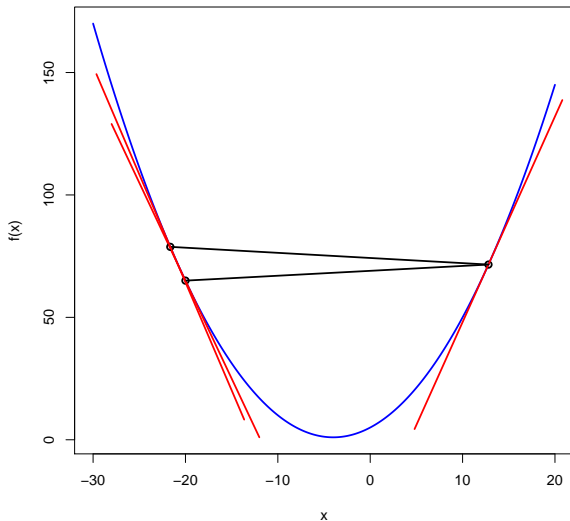
Logistic regression

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



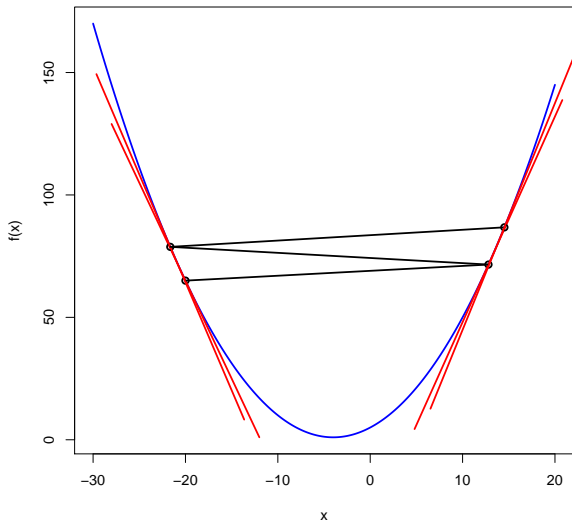
Logistic regression

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



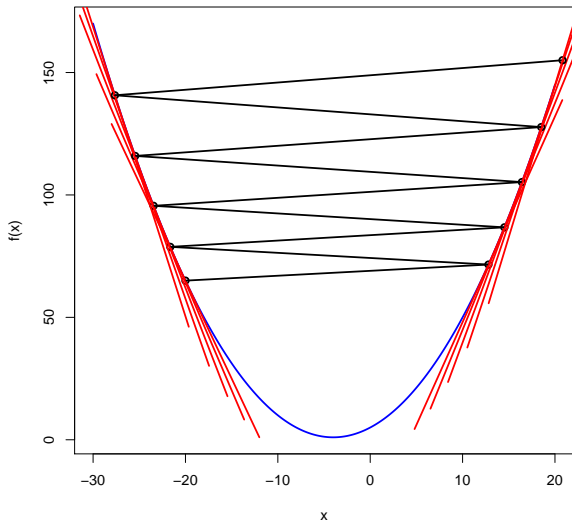
Logistic regression

- 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:

$$\frac{\partial \hat{L}_i}{\partial w_j} = -\beta_i y_i x_{ij},$$

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  $n$  {
         $w := w - \alpha \frac{\partial \hat{L}_i(w)}{\partial w}$ 
    }
}
```

Binary classification by boosting

- Iteratively learning base classifier on the data reweighted according to the error of previous base classifiers.

² Y. Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997
J. Friedman, T. Hastie, and R. Tibshirani. Additive Logistic Regression: a Statistical View of Boosting. *The Annals of Statistics*, 38(2), 2000

Binary classification by boosting

- Iteratively learning base classifier on the data reweighted according to the error of previous base classifiers.
- **AdaBoost**²
 - 1 $F_0 := 0$.
 - 2 For $k = 1, 2, \dots, K$:
$$f_k = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n e^{-y_i (F_{k-1}(\mathbf{x}_i) + f(\mathbf{x}_i))},$$
$$F_k := F_{k-1} + f_k.$$
 - 3 Return F_K .

² Y. Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997
J. Friedman, T. Hastie, and R. Tibshirani. Additive Logistic Regression: a Statistical View of Boosting. *The Annals of Statistics*, 38(2), 2000

Binary classification by boosting

- Effectively, we solve by the **coordinate descent method**:

$$\hat{f} = \arg \min_{f \in \text{lin}(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n e^{-y_i f(\mathbf{x}_i)}$$

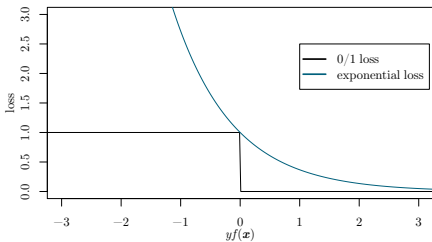
Binary classification by boosting

- Effectively, we solve by the **coordinate descent method**:

$$\hat{f} = \arg \min_{f \in \text{lin}(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n e^{-y_i f(\mathbf{x}_i)}$$

- We replaced 0/1 loss with **exponential loss**:

$$\ell_{\text{exp}}(y, f(\mathbf{x})) = e^{-yf(\mathbf{x})}.$$



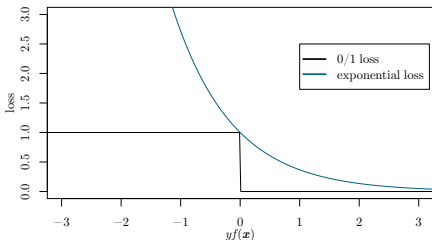
Binary classification by boosting

- Effectively, we solve by the **coordinate descent method**:

$$\hat{f} = \arg \min_{f \in \text{lin}(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n e^{-y_i f(\mathbf{x}_i)}$$

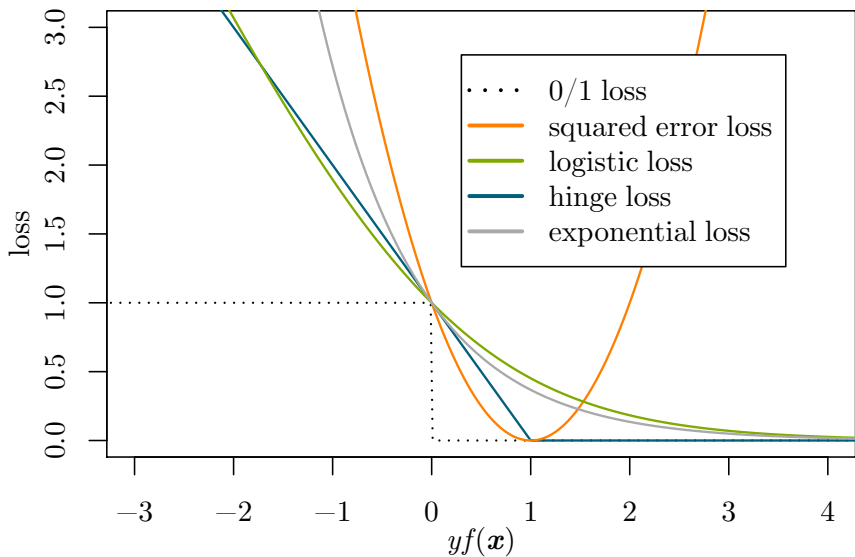
- We replaced 0/1 loss with **exponential loss**:

$$\ell_{\text{exp}}(y, f(\mathbf{x})) = e^{-yf(\mathbf{x})}.$$



- Once called “best off-the-shelf method for classification.” (Breiman)

Surrogate convex losses



Outline

- 1 Linear models for classification
- 2 Some statistical decision theory**
- 3 Classification calibrated losses
- 4 Generalization error

Statistical learning framework

Statistical learning framework

- **Input** $\mathbf{x} \in \mathcal{X}$ drawn from a distribution $P(\mathbf{x})$.
 - ▶ usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.

Statistical learning framework

- **Input** $\mathbf{x} \in \mathcal{X}$ drawn from a distribution $P(\mathbf{x})$.
 - ▶ usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- **Outcome** $y \in \mathcal{Y}$ drawn from a distribution $P(y | \mathbf{x})$.
 - ▶ target of our prediction: class label, real value, label vector, etc.,
 - ▶ alternative view: **example** (\mathbf{x}, y) drawn from $P(\mathbf{x}, y)$.

Statistical learning framework

- **Input** $\mathbf{x} \in \mathcal{X}$ drawn from a distribution $P(\mathbf{x})$.
 - ▶ usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- **Outcome** $y \in \mathcal{Y}$ drawn from a distribution $P(y | \mathbf{x})$.
 - ▶ target of our prediction: class label, real value, label vector, etc.,
 - ▶ alternative view: **example** (\mathbf{x}, y) drawn from $P(\mathbf{x}, y)$.
- **Prediction** $\hat{y} = h(\mathbf{x})$ by means of **prediction function** $h: \mathcal{X} \rightarrow \mathcal{Y}$.
 - ▶ h returns prediction $\hat{y} = h(\mathbf{x})$ for every input \mathbf{x} .

Statistical learning framework

- **Input** $\mathbf{x} \in \mathcal{X}$ drawn from a distribution $P(\mathbf{x})$.
 - ▶ usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- **Outcome** $y \in \mathcal{Y}$ drawn from a distribution $P(y | \mathbf{x})$.
 - ▶ target of our prediction: class label, real value, label vector, etc.,
 - ▶ alternative view: **example** (\mathbf{x}, y) drawn from $P(\mathbf{x}, y)$.
- **Prediction** $\hat{y} = h(\mathbf{x})$ by means of **prediction function** $h: \mathcal{X} \rightarrow \mathcal{Y}$.
 - ▶ h returns prediction $\hat{y} = h(\mathbf{x})$ for every input \mathbf{x} .
- **Loss** of our prediction: $\ell(y, \hat{y})$.
 - ▶ $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is a problem-specific **loss function**.

Statistical learning framework

- **Input** $\mathbf{x} \in \mathcal{X}$ drawn from a distribution $P(\mathbf{x})$.
 - ▶ usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- **Outcome** $y \in \mathcal{Y}$ drawn from a distribution $P(y | \mathbf{x})$.
 - ▶ target of our prediction: class label, real value, label vector, etc.,
 - ▶ alternative view: **example** (\mathbf{x}, y) drawn from $P(\mathbf{x}, y)$.
- **Prediction** $\hat{y} = h(\mathbf{x})$ by means of **prediction function** $h: \mathcal{X} \rightarrow \mathcal{Y}$.
 - ▶ h returns prediction $\hat{y} = h(\mathbf{x})$ for every input \mathbf{x} .
- **Loss** of our prediction: $\ell(y, \hat{y})$.
 - ▶ $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ is a problem-specific **loss function**.
- **Goal**: find a prediction function with small loss.

Risk

- **Goal:** minimize the **expected** loss over all examples (**risk**):

$$L_{\ell}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))].$$

Risk

- **Goal:** minimize the **expected** loss over all examples (**risk**):

$$L_\ell(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))].$$

- The **optimal** prediction function over all possible functions:

$$h^* = \arg \min_h L_\ell(h),$$

(so called **Bayes prediction function**).

Risk

- **Goal:** minimize the **expected** loss over all examples (**risk**):

$$L_\ell(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))].$$

- The **optimal** prediction function over all possible functions:

$$h^* = \arg \min_h L_\ell(h),$$

(so called **Bayes prediction function**).

- The smallest achievable risk (**Bayes risk**):

$$L_\ell^* = L_\ell(h^*).$$

Decomposition of risk

$$L_\ell(h) = \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))]$$

Decomposition of risk

$$\begin{aligned}L_{\ell}(h) &= \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))] \\ &= \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(\mathbf{x})) P(\mathbf{x}, y) d\mathbf{x} dy\end{aligned}$$

Decomposition of risk

$$\begin{aligned}L_{\ell}(h) &= \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))] \\&= \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(\mathbf{x})) P(\mathbf{x}, y) d\mathbf{x} dy \\&= \int_{\mathcal{X}} \left(\int_{\mathcal{Y}} \ell(y, h(\mathbf{x})) P(y | \mathbf{x}) dy \right) P(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Decomposition of risk

$$\begin{aligned}L_{\ell}(h) &= \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))] \\&= \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(\mathbf{x})) P(\mathbf{x}, y) d\mathbf{x} dy \\&= \int_{\mathcal{X}} \left(\int_{\mathcal{Y}} \ell(y, h(\mathbf{x})) P(y | \mathbf{x}) dy \right) P(\mathbf{x}) d\mathbf{x} \\&= \mathbb{E}_{\mathbf{x}} [L_{\ell}(h | \mathbf{x})].\end{aligned}$$

Decomposition of risk

$$\begin{aligned}L_{\ell}(h) &= \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))] \\&= \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(\mathbf{x})) P(\mathbf{x}, y) d\mathbf{x} dy \\&= \int_{\mathcal{X}} \left(\int_{\mathcal{Y}} \ell(y, h(\mathbf{x})) P(y | \mathbf{x}) dy \right) P(\mathbf{x}) d\mathbf{x} \\&= \mathbb{E}_{\mathbf{x}} [L_{\ell}(h | \mathbf{x})].\end{aligned}$$

- $L_{\ell}(h | \mathbf{x})$ is the **conditional risk** of $\hat{y} = h(\mathbf{x})$ at \mathbf{x} .

Decomposition of risk

$$\begin{aligned}L_{\ell}(h) &= \mathbb{E}_{(\mathbf{x}, y)} [\ell(y, h(\mathbf{x}))] \\&= \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, h(\mathbf{x})) P(\mathbf{x}, y) d\mathbf{x} dy \\&= \int_{\mathcal{X}} \left(\int_{\mathcal{Y}} \ell(y, h(\mathbf{x})) P(y | \mathbf{x}) dy \right) P(\mathbf{x}) d\mathbf{x} \\&= \mathbb{E}_{\mathbf{x}} [L_{\ell}(h | \mathbf{x})].\end{aligned}$$

- $L_{\ell}(h | \mathbf{x})$ is the **conditional risk** of $\hat{y} = h(\mathbf{x})$ at \mathbf{x} .
- Bayes prediction **minimizes the conditional risk** for every \mathbf{x} :

$$h^*(\mathbf{x}) = \arg \min_h L_{\ell}(h | \mathbf{x}).$$

Example — binary classification and 0/1 loss

- The outcome y is **binary**, $y \in \{-1, 1\}$.

Example — binary classification and 0/1 loss

- The outcome y is **binary**, $y \in \{-1, 1\}$.
- **0/1** loss function:

$$\ell_{0/1}(y, \hat{y}) = \begin{cases} 0, & \text{if } y = \hat{y}, \\ 1, & \text{otherwise.} \end{cases}$$

Example — binary classification and 0/1 loss

- The outcome y is **binary**, $y \in \{-1, 1\}$.
- **0/1** loss function:

$$\ell_{0/1}(y, \hat{y}) = \begin{cases} 0, & \text{if } y = \hat{y}, \\ 1, & \text{otherwise.} \end{cases}$$

- Define $\eta(\mathbf{x}) = P(y = 1|\mathbf{x})$. The conditional 0/1 risk at \mathbf{x} is:

$$L_{0/1}(h|\mathbf{x}) = \eta(\mathbf{x})\mathbb{I}[h(\mathbf{x}) = -1] + (1 - \eta(\mathbf{x}))\mathbb{I}[h(\mathbf{x}) = 1].$$

Example — binary classification and 0/1 loss

- The outcome y is **binary**, $y \in \{-1, 1\}$.
- **0/1** loss function:

$$\ell_{0/1}(y, \hat{y}) = \begin{cases} 0, & \text{if } y = \hat{y}, \\ 1, & \text{otherwise.} \end{cases}$$

- Define $\eta(\mathbf{x}) = P(y = 1|\mathbf{x})$. The conditional 0/1 risk at \mathbf{x} is:

$$L_{0/1}(h|\mathbf{x}) = \eta(\mathbf{x})\mathbb{I}[h(\mathbf{x}) = -1] + (1 - \eta(\mathbf{x}))\mathbb{I}[h(\mathbf{x}) = 1].$$

- The **Bayes classifier**:

$$h^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \eta(\mathbf{x}) > 1 - \eta(\mathbf{x}) \\ -1 & \text{if } \eta(\mathbf{x}) < 1 - \eta(\mathbf{x}) \end{cases} = \text{sgn}(\eta(\mathbf{x}) - 1/2),$$

and the **Bayes conditional risk**:

$$L_{\ell}(h^* | \mathbf{x}) = \min\{\eta(\mathbf{x}), 1 - \eta(\mathbf{x})\}.$$

Regret

- Bayes risk measures the **difficulty** of the problems.
- Difficult data \implies high Bayes risk \implies every h has high risk.

Regret

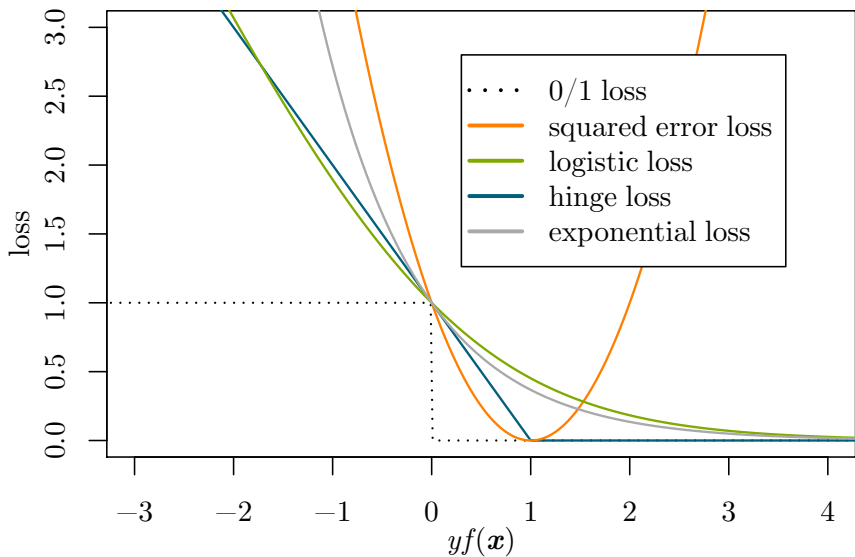
- Bayes risk measures the **difficulty** of the problems.
- Difficult data \implies high Bayes risk \implies every h has high risk.
- It is not the risk what captures the prediction accuracy of h , but rather **how close the risk of h is to the risk of h^*** .

Regret

- Bayes risk measures the **difficulty** of the problems.
- Difficult data \implies high Bayes risk \implies every h has high risk.
- It is not the risk what captures the prediction accuracy of h , but rather **how close the risk of h is to the risk of h^*** .
- The **regret (excess risk)** of a prediction function h :

$$\text{Reg}_\ell(h) = L_\ell(h) - L_\ell^*$$

Surrogate convex losses

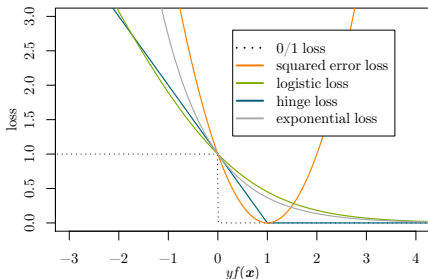


Surrogate convex losses

In each case, $\ell(y, f(\mathbf{x}))$ is a function of a **single** variable:

$$\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x})).$$

We call $yf(\mathbf{x})$ the **margin**, and such ℓ – **margin-based loss function**.

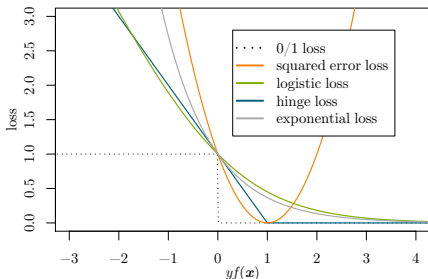


Surrogate convex losses

In each case, $\ell(y, f(\mathbf{x}))$ is a function of a **single** variable:

$$\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x})).$$

We call $yf(\mathbf{x})$ the **margin**, and such ℓ – **margin-based loss function**.



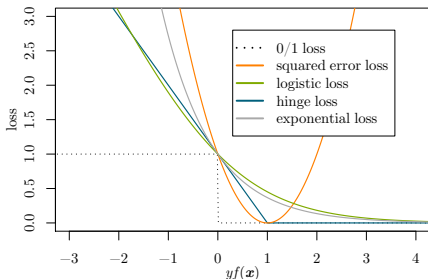
- All losses are **convex** and can be scaled to **upperbound** 0/1 loss.

Surrogate convex losses

In each case, $\ell(y, f(\mathbf{x}))$ is a function of a **single** variable:

$$\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x})).$$

We call $yf(\mathbf{x})$ the **margin**, and such ℓ – **margin-based loss function**.



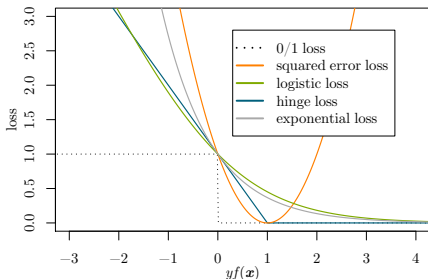
- All losses are **convex** and can be scaled to **upperbound** 0/1 loss.
- Minimizing such an upper bound gives a bound on **0/1 risk**, but **not** on the **0/1 regret**!

Surrogate convex losses

In each case, $\ell(y, f(\mathbf{x}))$ is a function of a **single** variable:

$$\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x})).$$

We call $yf(\mathbf{x})$ the **margin**, and such ℓ – **margin-based loss function**.



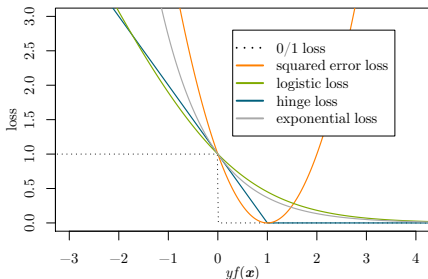
- All losses are **convex** and can be scaled to **upperbound** 0/1 loss.
- Minimizing such an upper bound gives a bound on **0/1 risk**, but **not** on the **0/1 regret**!
- Do we actually solve the original classification problem by minimizing the surrogate loss?

Surrogate convex losses

In each case, $\ell(y, f(\mathbf{x}))$ is a function of a **single** variable:

$$\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x})).$$

We call $yf(\mathbf{x})$ the **margin**, and such ℓ – **margin-based loss function**.



- All losses are **convex** and can be scaled to **upperbound** 0/1 loss.
- Minimizing such an upper bound gives a bound on **0/1 risk**, but **not** on the **0/1 regret**!
- Do we actually solve the original classification problem by minimizing the surrogate loss?
- Need a **theory** for classification with **surrogate convex losses**.

Outline

- 1 Linear models for classification
- 2 Some statistical decision theory
- 3 Classification calibrated losses**
- 4 Generalization error

Calibration

- **Task loss** ℓ_T , the one to be actually minimized.
 - ▶ e.g., 0/1 loss in binary classification.

Calibration

- **Task loss** ℓ_T , the one to be actually minimized.
 - ▶ e.g., 0/1 loss in binary classification.
- **Surrogate (proxy) loss** ℓ_S , the one which we use in learning.
 - ▶ e.g., logistic loss in binary classification.

Calibration

- **Task loss** ℓ_T , the one to be actually minimized.
 - ▶ e.g., 0/1 loss in binary classification.
- **Surrogate (proxy) loss** ℓ_S , the one which we use in learning.
 - ▶ e.g., logistic loss in binary classification.
- We say a surrogate loss ℓ_S is **calibrated (consistent)** for a task loss ℓ_T if for any sequence of functions $\{f_n\}$, and any distribution $P(\mathbf{x}, y)$:

$$L_{\ell_S}(f_n) \rightarrow \min_f L_{\ell_S}(f) \implies L_{\ell_T}(f_n) \rightarrow \min_f L_{\ell_T}(f).$$

or equivalently:

$$\text{Reg}_{\ell_S}(f_n) \rightarrow 0 \implies \text{Reg}_{\ell_T}(f_n) \rightarrow 0.$$

Calibration

- **Task loss** ℓ_T , the one to be actually minimized.
 - ▶ e.g., 0/1 loss in binary classification.
- **Surrogate (proxy) loss** ℓ_S , the one which we use in learning.
 - ▶ e.g., logistic loss in binary classification.
- We say a surrogate loss ℓ_S is **calibrated (consistent)** for a task loss ℓ_T if for any sequence of functions $\{f_n\}$, and any distribution $P(\mathbf{x}, y)$:

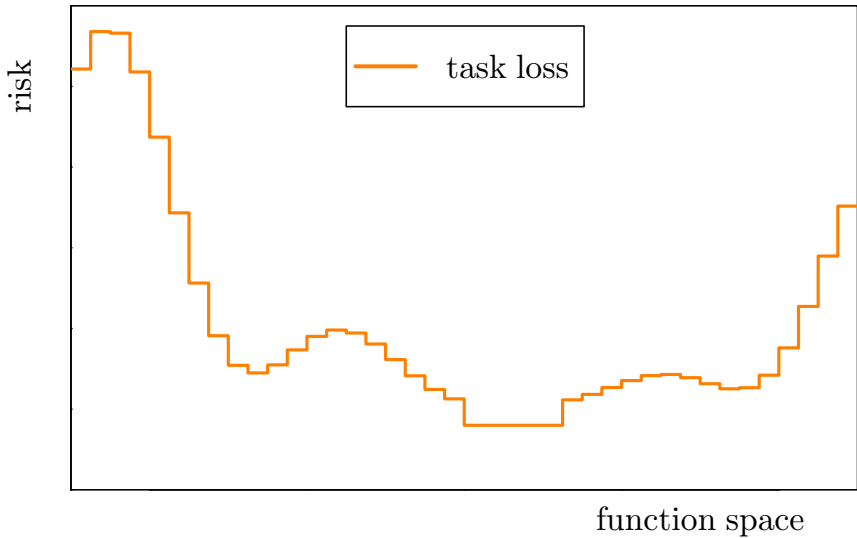
$$L_{\ell_S}(f_n) \rightarrow \min_f L_{\ell_S}(f) \implies L_{\ell_T}(f_n) \rightarrow \min_f L_{\ell_T}(f).$$

or equivalently:

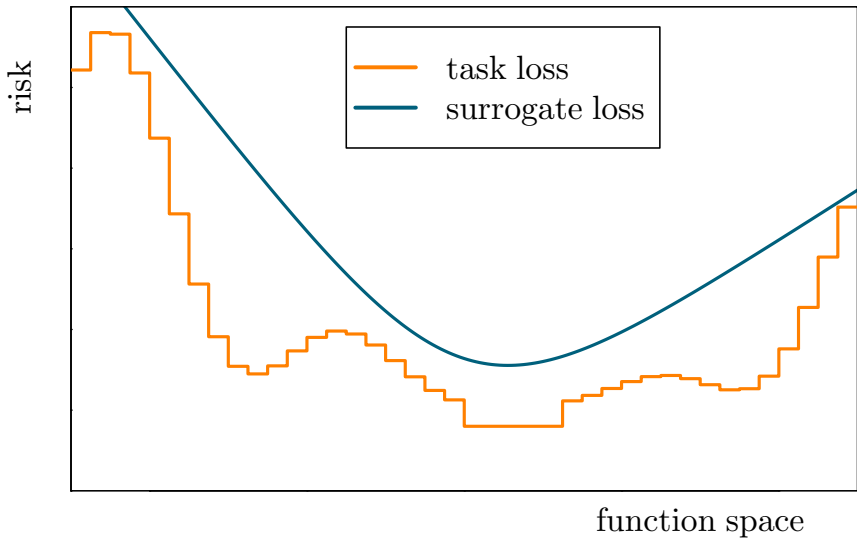
$$\text{Reg}_{\ell_S}(f_n) \rightarrow 0 \implies \text{Reg}_{\ell_T}(f_n) \rightarrow 0.$$

- Informally: if we converge to the prediction function **optimal with respect to ℓ_S** , then we also converge to the prediction function **optimal with respect to ℓ_T** .

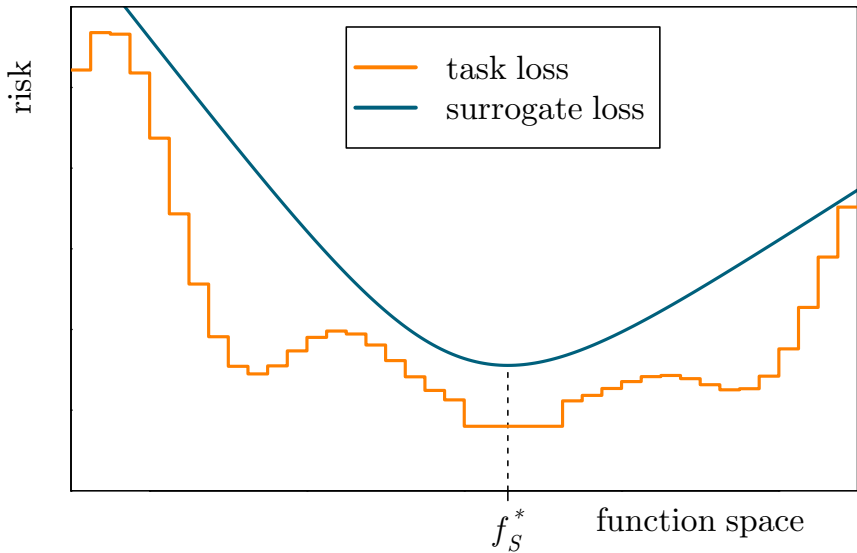
Calibrated loss



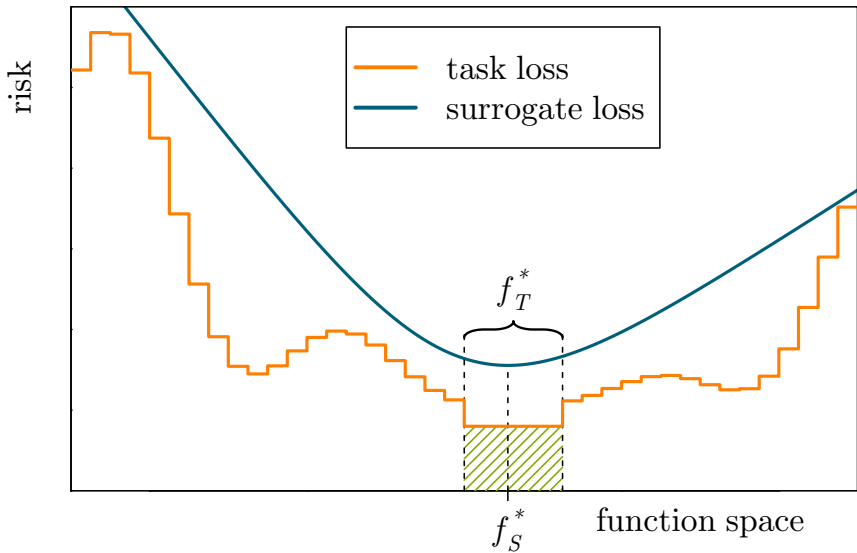
Calibrated loss



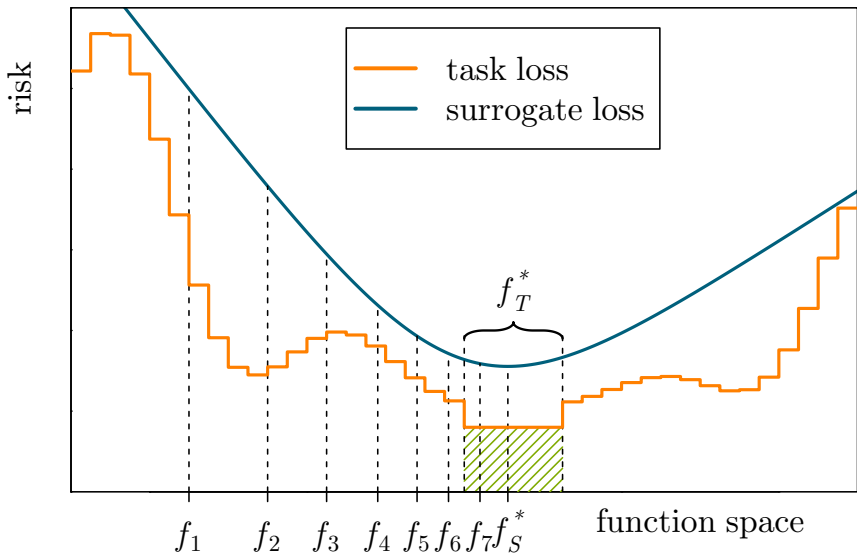
Calibrated loss



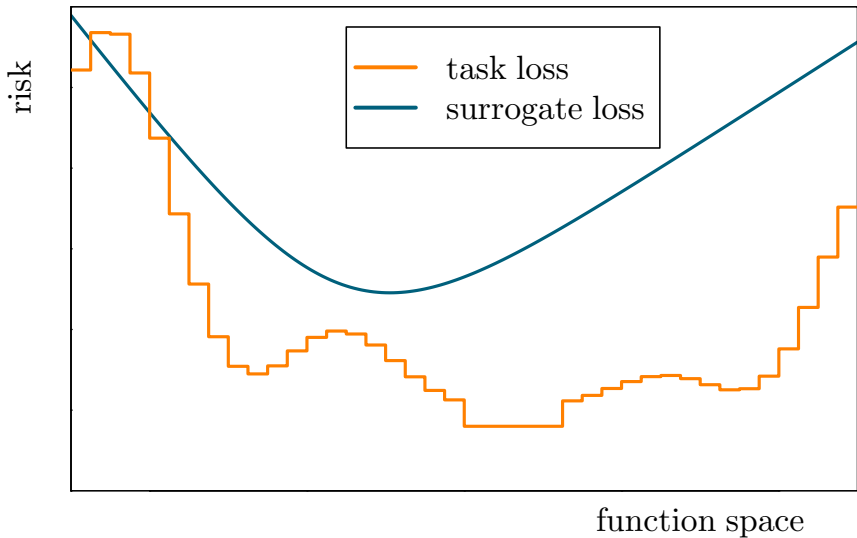
Calibrated loss



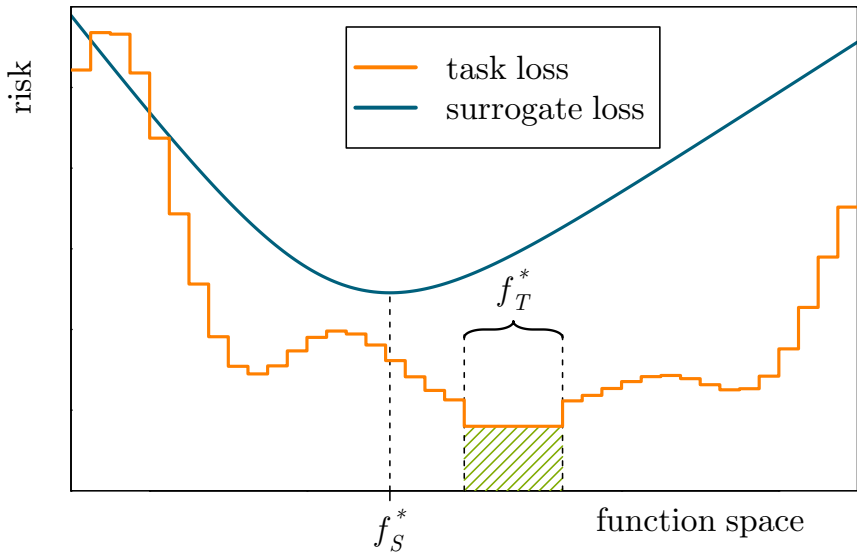
Calibrated loss



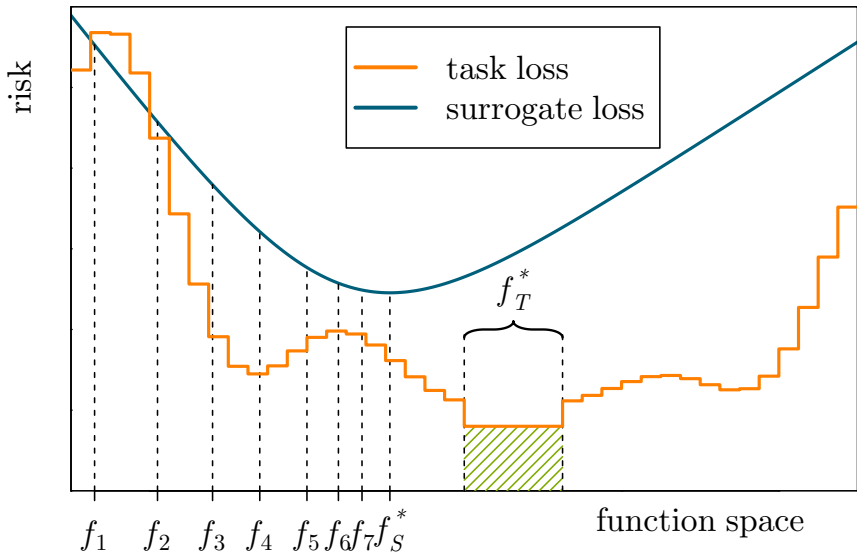
Non-calibrated loss



Non-calibrated loss



Non-calibrated loss



Classification calibration

- We say a surrogate loss ℓ is **classification calibrated (consistent)** if for any sequence of functions $\{f_n\}$ and any distribution $P(x, y)$:

$$\text{Reg}_\ell(f_n) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n) \rightarrow 0.$$

Classification calibration

- We say a surrogate loss ℓ is **classification calibrated (consistent)** if for any sequence of functions $\{f_n\}$ and any distribution $P(\mathbf{x}, y)$:

$$\text{Reg}_\ell(f_n) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n) \rightarrow 0.$$

- We only consider **margin-based losses**: $\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x}))$.

Classification calibration

- We say a surrogate loss ℓ is **classification calibrated (consistent)** if for any sequence of functions $\{f_n\}$ and any distribution $P(\mathbf{x}, y)$:

$$\text{Reg}_\ell(f_n) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n) \rightarrow 0.$$

- We only consider **margin-based losses**: $\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x}))$.
- Holds for any $P(\mathbf{x}, y) \implies$ must hold conditionally for any \mathbf{x} :

Question: why?

$$\begin{aligned} \text{Reg}_\ell(f_n|\mathbf{x}) \rightarrow 0 &\implies \text{Reg}_{0/1}(f_n|\mathbf{x}) \rightarrow 0, && \text{or} \\ L_\ell(f_n|\mathbf{x}) \rightarrow \min_f L_\ell(f|\mathbf{x}) &\implies L_{0/1}(f_n|\mathbf{x}) \rightarrow \min_f L_{0/1}(f|\mathbf{x}). \end{aligned}$$

Classification calibration

- We say a surrogate loss ℓ is **classification calibrated (consistent)** if for any sequence of functions $\{f_n\}$ and any distribution $P(\mathbf{x}, y)$:

$$\text{Reg}_\ell(f_n) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n) \rightarrow 0.$$

- We only consider **margin-based losses**: $\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x}))$.
- Holds for any $P(\mathbf{x}, y) \implies$ must hold conditionally for any \mathbf{x} :

Question: why?

$$\text{Reg}_\ell(f_n|\mathbf{x}) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n|\mathbf{x}) \rightarrow 0, \quad \text{or}$$

$$L_\ell(f_n|\mathbf{x}) \rightarrow \min_f L_\ell(f|\mathbf{x}) \implies L_{0/1}(f_n|\mathbf{x}) \rightarrow \min_f L_{0/1}(f|\mathbf{x}).$$

- **Abstraction**: $L_\ell(f|\mathbf{x})$ depends **only on two numbers**: $f(\mathbf{x})$ and $\eta(\mathbf{x})$. Also, $\min_f L_\ell(f|\mathbf{x})$ only depends on a **single number** $\eta(\mathbf{x})$.

Classification calibration

- We say a surrogate loss ℓ is **classification calibrated (consistent)** if for any sequence of functions $\{f_n\}$ and any distribution $P(\mathbf{x}, y)$:

$$\text{Reg}_\ell(f_n) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n) \rightarrow 0.$$

- We only consider **margin-based losses**: $\ell(y, f(\mathbf{x})) = \ell(yf(\mathbf{x}))$.
- Holds for any $P(\mathbf{x}, y) \implies$ must hold conditionally for any \mathbf{x} :

Question: why?

$$\text{Reg}_\ell(f_n|\mathbf{x}) \rightarrow 0 \implies \text{Reg}_{0/1}(f_n|\mathbf{x}) \rightarrow 0, \quad \text{or}$$

$$L_\ell(f_n|\mathbf{x}) \rightarrow \min_f L_\ell(f|\mathbf{x}) \implies L_{0/1}(f_n|\mathbf{x}) \rightarrow \min_f L_{0/1}(f|\mathbf{x}).$$

- **Abstraction**: $L_\ell(f|\mathbf{x})$ depends **only on two numbers**: $f(\mathbf{x})$ and $\eta(\mathbf{x})$. Also, $\min_f L_\ell(f|\mathbf{x})$ only depends on a **single number** $\eta(\mathbf{x})$.
- We drop the dependence on \mathbf{x} and test convergence for any sequence of **numbers** $\{f_n\}$ and any **number** η .

Classification calibration³

- A shorthand notation for **conditional 0/1-risk** $L_{0/1}(f|\mathbf{x})$:

$$C_{\eta}^{0/1}(f) = \eta \llbracket f < 0 \rrbracket + (1 - \eta) \llbracket f > 0 \rrbracket + \frac{1}{2} \llbracket f = 0 \rrbracket.$$

³ P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006

Classification calibration³

- A shorthand notation for **conditional 0/1-risk** $L_{0/1}(f|\mathbf{x})$:

$$C_{\eta}^{0/1}(f) = \eta \llbracket f < 0 \rrbracket + (1 - \eta) \llbracket f > 0 \rrbracket + \frac{1}{2} \llbracket f = 0 \rrbracket.$$

- $C_{\eta}^{0/1}(f)$ is **minimized** by any f^* , s.t. $\text{sgn}(f^*) = \text{sgn}(\eta - 1/2)$ for $\eta \neq 1/2$, and f^* **arbitrary** for $\eta = 1/2$.
 - ▶ The **Bayes (0/1) classifier** f^* is **non unique!**

³ P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006

Classification calibration³

- A shorthand notation for **conditional 0/1-risk** $L_{0/1}(f|\mathbf{x})$:

$$C_{\eta}^{0/1}(f) = \eta \llbracket f < 0 \rrbracket + (1 - \eta) \llbracket f > 0 \rrbracket + \frac{1}{2} \llbracket f = 0 \rrbracket.$$

- $C_{\eta}^{0/1}(f)$ is **minimized** by any f^* , s.t. $\text{sgn}(f^*) = \text{sgn}(\eta - 1/2)$ for $\eta \neq 1/2$, and f^* **arbitrary** for $\eta = 1/2$.
 - ▶ The **Bayes (0/1) classifier** f^* is **non unique!**
- The **Bayes 0/1-risk**:

$$\min_f C_{\eta}^{0/1}(f) = \min\{\eta, 1 - \eta\},$$

³ P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006

Classification calibration

- A shorthand notation for **conditional ℓ -risk** $L_\ell(f|\mathbf{x})$:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f).$$

Classification calibration

- A shorthand notation for **conditional ℓ -risk** $L_\ell(f|\mathbf{x})$:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f).$$

- **Bayes conditional ℓ -risk**:

$$H(\eta) := \min_f C_\eta(f).$$

Classification calibration

- A shorthand notation for **conditional ℓ -risk** $L_\ell(f|\mathbf{x})$:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f).$$

- **Bayes conditional ℓ -risk**:

$$H(\eta) := \min_f C_\eta(f).$$

- Loss ℓ is **classification calibrated** if and only if for any $\{f_n\}$ and η :

$$C_\eta(f_n) \rightarrow H(\eta) \implies C_\eta^{0/1}(f_n) \rightarrow \min\{\eta, 1 - \eta\}.$$

Classification calibration

- To have $C_\eta^{0/1}(f_n) \rightarrow \min\{\eta, 1 - \eta\}$, for any $\eta \neq 1/2$, f_n must **eventually** pick the **correct sign**, i.e. for large n ,

$$\text{sgn}(f_n) = \text{sgn}(\eta - 1/2), \quad \text{or} \quad f_n(\eta - 1/2) > 0.$$

Classification calibration

- To have $C_\eta^{0/1}(f_n) \rightarrow \min\{\eta, 1 - \eta\}$, for any $\eta \neq 1/2$, f_n must **eventually** pick the **correct sign**, i.e. for large n ,

$$\text{sgn}(f_n) = \text{sgn}(\eta - 1/2), \quad \text{or} \quad f_n(\eta - 1/2) > 0.$$

- Since $C_\eta(f_n) \rightarrow \min_f C_\eta(f)$, this means any minimizer f^* of $C_\eta(f)$ must satisfy $f^*(\eta - 1/2) > 0$.

Classification calibration

- To have $C_\eta^{0/1}(f_n) \rightarrow \min\{\eta, 1 - \eta\}$, for any $\eta \neq 1/2$, f_n must **eventually** pick the **correct sign**, i.e. for large n ,

$$\text{sgn}(f_n) = \text{sgn}(\eta - 1/2), \quad \text{or} \quad f_n(\eta - 1/2) > 0.$$

- Since $C_\eta(f_n) \rightarrow \min_f C_\eta(f)$, this means any minimizer f^* of $C_\eta(f)$ must satisfy $f^*(\eta - 1/2) > 0$, i.e.:

$$H(\eta) = \min_f C_\eta(f) < \min_{f: f(\eta-1/2) \leq 0} C_\eta(f) =: H^-(\eta).$$

Classification calibration

- To have $C_\eta^{0/1}(f_n) \rightarrow \min\{\eta, 1 - \eta\}$, for any $\eta \neq 1/2$, f_n must **eventually** pick the **correct sign**, i.e. for large n ,

$$\text{sgn}(f_n) = \text{sgn}(\eta - 1/2), \quad \text{or} \quad f_n(\eta - 1/2) > 0.$$

- Since $C_\eta(f_n) \rightarrow \min_f C_\eta(f)$, this means any minimizer f^* of $C_\eta(f)$ must satisfy $f^*(\eta - 1/2) > 0$, i.e.:

$$H(\eta) = \min_f C_\eta(f) < \min_{f: f(\eta-1/2) \leq 0} C_\eta(f) =: H^-(\eta).$$

- **Fact:** ℓ is **classification calibrated if and only if** for any $\eta \neq 1/2$,

$$H(\eta) < H^-(\eta).$$

The minimum of the conditional ℓ -risk is always achieved on the **correct side of 0**.

Classification calibration

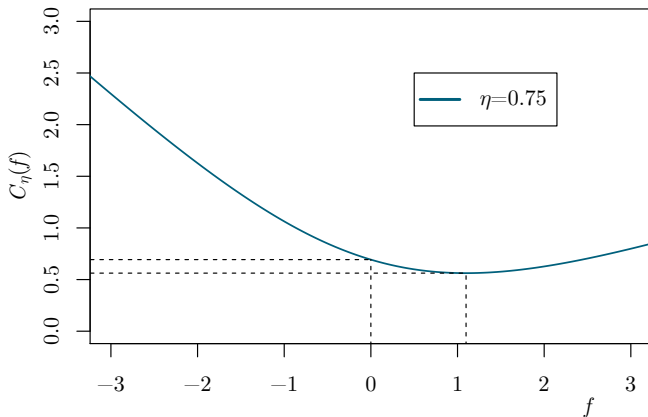
Loss ℓ is **classification calibrated** if and only if $H(\eta) < H^-(\eta)$, i.e.

$$\min_f C_\eta(f) < \min_{f: f(\eta-1/2) \leq 0} C_\eta(f).$$

Classification calibration

Loss ℓ is **classification calibrated** if and only if $H(\eta) < H^-(\eta)$, i.e.

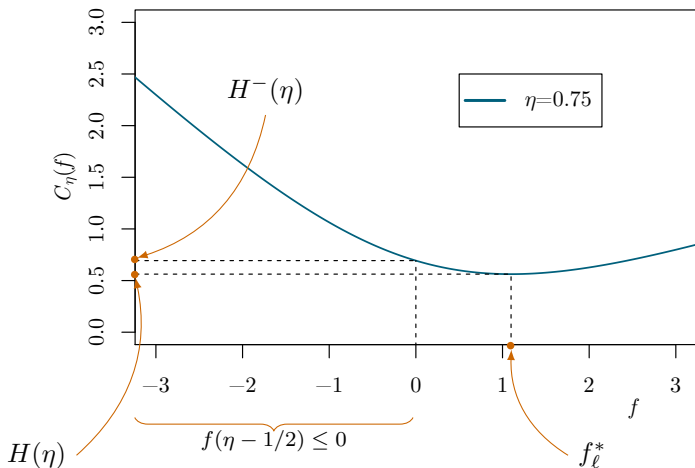
$$\min_f C_\eta(f) < \min_{f: f(\eta-1/2) \leq 0} C_\eta(f).$$



Classification calibration

Loss ℓ is **classification calibrated** if and only if $H(\eta) < H^-(\eta)$, i.e.

$$\min_f C_\eta(f) < \min_{f: f(\eta-1/2) \leq 0} C_\eta(f).$$



Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



First take $\eta > 1/2$ and assume $\ell'(0) < 0$.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



First take $\eta > 1/2$ and assume $\ell'(0) < 0$.

$\ell(f)$ is convex, then so is $C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f)$ (why?) .

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



First take $\eta > 1/2$ and assume $\ell'(0) < 0$.

$\ell(f)$ is convex, then so is $C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f)$ (why?) .

Moreover,

$$C'_\eta(0) = (2\eta - 1)\ell'(0) < 0.$$

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



First take $\eta > 1/2$ and assume $\ell'(0) < 0$.

$\ell(f)$ is convex, then so is $C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f)$ (why?).

Moreover,

$$C'_\eta(0) = (2\eta - 1)\ell'(0) < 0.$$

For convex function the derivative is **increasing**, so the minimum must be on the **positive side** of 0.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)



First take $\eta > 1/2$ and assume $\ell'(0) < 0$.

$\ell(f)$ is convex, then so is $C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f)$ (why?).

Moreover,

$$C'_\eta(0) = (2\eta - 1)\ell'(0) < 0.$$

For convex function the derivative is **increasing**, so the minimum must be on the **positive side** of 0.

Similar argument for $\eta < 1/2$.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

\implies

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

\implies

First take $\eta > 1/2$ and assume $\ell(f)$ is **calibrated**.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

\implies

First take $\eta > 1/2$ and assume $\ell(f)$ is **calibrated**.

Convexity of ℓ implies that for any f , $\ell(f) \geq \ell(0) + \ell'(0)f$ (tangent to the function chart at 0 is always below the chart).

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

\implies

First take $\eta > 1/2$ and assume $\ell(f)$ is **calibrated**.

Convexity of ℓ implies that for any f , $\ell(f) \geq \ell(0) + \ell'(0)f$ (tangent to the function chart at 0 is always below the chart).

We get:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f) \geq \ell(0) + \ell'(0)(2\eta - 1).$$

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

⇒

First take $\eta > 1/2$ and assume $\ell(f)$ is **calibrated**.

Convexity of ℓ implies that for any f , $\ell(f) \geq \ell(0) + \ell'(0)f$ (tangent to the function chart at 0 is always below the chart).

We get:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f) \geq \ell(0) + \ell'(0)(2\eta - 1).$$

The r.h.s. must be **less than** $\ell(0)$ (so that $\ell'(0) < 0$). Otherwise, $C_\eta(f) \geq \ell(0) = C_\eta(0)$, which implies $f = 0$ is the minimizer of ℓ , so $\ell(f)$ is not calibrated.

Calibration — simplified condition

A convex loss $\ell(f)$ is classification calibrated **if and only if** it is **differentiable at 0** and $\ell'(0) < 0$.

Proof (sketch):

(only show the proof for ℓ differentiable at 0)

⇒

First take $\eta > 1/2$ and assume $\ell(f)$ is **calibrated**.

Convexity of ℓ implies that for any f , $\ell(f) \geq \ell(0) + \ell'(0)f$ (tangent to the function chart at 0 is always below the chart).

We get:

$$C_\eta(f) = \eta\ell(f) + (1 - \eta)\ell(-f) \geq \ell(0) + \ell'(0)(2\eta - 1).$$

The r.h.s. must be **less than** $\ell(0)$ (so that $\ell'(0) < 0$). Otherwise, $C_\eta(f) \geq \ell(0) = C_\eta(0)$, which implies $f = 0$ is the minimizer of ℓ , so $\ell(f)$ is not calibrated. Similar argument for $\eta < 1/2$.

Examples

Question

Prove that all the following losses are calibrated and **derive** their Bayes classifiers:

- squared error loss $\ell(f) = (1 - f)^2$,
- logistic loss $\ell(f) = \log(1 + e^{-f})$,
- hinge loss $\ell(f) = \max\{0, 1 - f\}$,
- exponential loss $\ell(f) = e^{-f}$.

Example — squared error loss

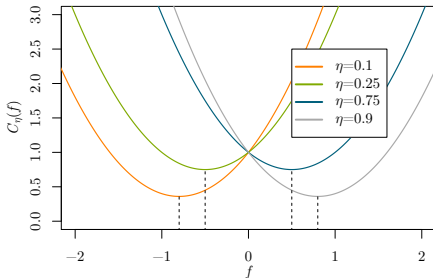
$$\ell(f) = (1 - f)^2$$

Conditional risk:

$$\begin{aligned}C_\eta(f) &= \eta(1 - f)^2 + (1 - \eta)(1 + f)^2 \\ &= 1 + f^2 + 2f(1 - 2\eta).\end{aligned}$$

Minimized by:

$$f^* = 2\eta - 1$$



Calibration:

$$\ell'(f) = -2(1 - f) \quad \implies \quad \ell'(0) = -2 < 0$$

$$H(\eta) = C_\eta(f^*) = 4\eta(1 - \eta),$$

$$H^-(\eta) = C_\eta(0) = 1 > H(\eta).$$

Example — logistic loss

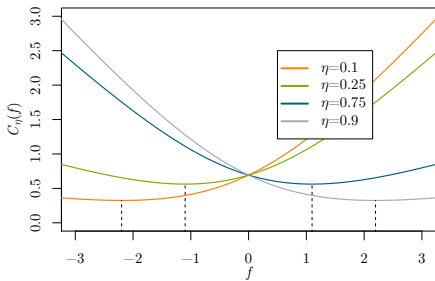
$$\ell(f) = \log(1 + e^{-f})$$

Conditional risk:

$$C_{\eta}(f) = \eta \log(1 + e^{-f}) + (1 - \eta) \log(1 + e^f).$$

Minimized by:

$$f^* = \log \frac{\eta}{1 - \eta}.$$



Calibration:

$$\ell'(f) = -\frac{1}{1 + e^f} \quad \implies \quad \ell'(0) = -1/2 < 0$$

$$H(\eta) = C_{\eta}(f^*) = -\eta \log \eta - (1 - \eta) \log(1 - \eta)$$

$$H^-(\eta) = C_{\eta}(0) = \log 2 > H(\eta).$$

Example — hinge loss

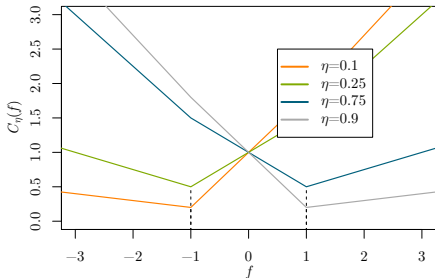
$$\ell(f) = \max\{0, 1 - f\}$$

Conditional risk:

$$C_\eta(f) = \eta \max\{0, 1 - f\} \\ + (1 - \eta) \max\{0, 1 + f\}.$$

Minimized by:

$$f^* = \operatorname{sgn}(\eta - 1/2)$$



Calibration:

$$\ell'(f) = -1 \llbracket f < 1 \rrbracket \quad \implies \quad \ell'(0) = -1 < 0$$

$$H(\eta) = C_\eta(f^*) = 2 \min\{\eta, 1 - \eta\}$$

$$H^-(\eta) = C_\eta(0) = 1 > H(\eta).$$

Example — exponential loss

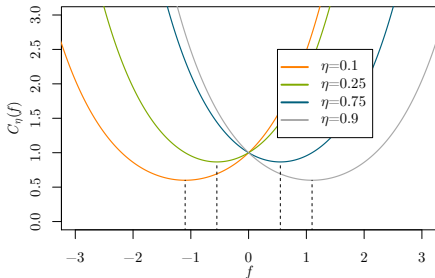
$$\ell(f) = e^{-f}$$

Conditional risk:

$$C_\eta(f) = \eta e^{-f} + (1 - \eta)e^f$$

Minimized by:

$$f^* = \frac{1}{2} \log \frac{\eta}{1 - \eta}$$



Calibration:

$$\ell'(f) = -e^{-f} \implies \ell'(0) = -1 < 0$$

$$H(\eta) = C_\eta(f^*) = 2\sqrt{\eta(1 - \eta)}$$

$$H^-(\eta) = C_\eta(0) = 1 > H(\eta).$$

Calibration of surrogate losses

- All considered surrogate losses are **classification calibrated**.

Calibration of surrogate losses

- All considered surrogate losses are **classification calibrated**.
- In each case f^* is an increasing function of η .

Calibration of surrogate losses

- All considered surrogate losses are **classification calibrated**.
- In each case f^* is an increasing function of η .
- In all cases except hinge loss, $f^*(\eta)$ is **invertible**
 \implies we can estimate η by inverting the relation.

Calibration of surrogate losses

- All considered surrogate losses are **classification calibrated**.
- In each case f^* is an increasing function of η .
- In all cases except hinge loss, $f^*(\eta)$ is **invertible**
 \implies we can estimate η by inverting the relation.

loss	$f^*(\eta)$	$\eta(f^*)$
squared error	$2\eta - 1$	$\frac{1+f^*}{2}$
logistic	$\log \frac{\eta}{1-\eta}$	$\frac{1}{1+e^{-f^*}}$
exponential	$\frac{1}{2} \log \frac{\eta}{1-\eta}$	$\frac{1}{1+e^{-2f^*}}$
hinge	$\text{sgn}(\eta - 1/2)$	doesn't exist

Calibration of surrogate losses

- All considered surrogate losses are **classification calibrated**.
- In each case f^* is an increasing function of η .
- In all cases except hinge loss, $f^*(\eta)$ is **invertible**
 \implies we can estimate η by inverting the relation.

loss	$f^*(\eta)$	$\eta(f^*)$
squared error	$2\eta - 1$	$\frac{1+f^*}{2}$
logistic	$\log \frac{\eta}{1-\eta}$	$\frac{1}{1+e^{-f^*}}$
exponential	$\frac{1}{2} \log \frac{\eta}{1-\eta}$	$\frac{1}{1+e^{-2f^*}}$
hinge	$\text{sgn}(\eta - 1/2)$	doesn't exist

- Minimizing hinge loss **does not** lead to probability estimation.

Outline

- 1 Linear models for classification
- 2 Some statistical decision theory
- 3 Classification calibrated losses
- 4 Generalization error**

Training vs. test error

Training error of function h :

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i))$$

\implies Average error on a sample of size n .

Training vs. test error

Training error of function h :

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i))$$

\implies Average error on a sample of size n .

Test error of function h :

$$L(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))]$$

\implies Expected error h makes on a randomly drawn instance.

Generalization error of function h : $L(h) - \widehat{L}(h)$
(test error – training error)

Training vs. test error

Training error of function h :

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i))$$

\implies Average error on a sample of size n .

Test error of function h :

$$L(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))]$$

\implies Expected error h makes on a randomly drawn instance.

Generalization error of function h : $L(h) - \widehat{L}(h)$
(test error – training error)

- For fixed h , $\widehat{L}(h) \simeq L(h)$ (law of large numbers).

Training vs. test error

Training error of function h :

$$\widehat{L}(h) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i))$$

\implies Average error on a sample of size n .

Test error of function h :

$$L(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y, h(\mathbf{x}))]$$

\implies Expected error h makes on a randomly drawn instance.

Generalization error of function h : $L(h) - \widehat{L}(h)$
(test error – training error)

- For fixed h , $\widehat{L}(h) \simeq L(h)$ (law of large numbers).
- But h is not fixed, it is **chosen** based on the training sample!

Empirical Risk Minimization (ERM)

- Choose a prediction function \hat{h} which minimizes the loss on the training data within some **restricted** class of functions \mathcal{H} .

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i)).$$

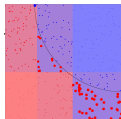
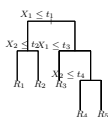
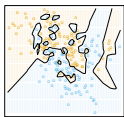
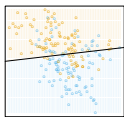
- The average loss on the training data is called **empirical risk** $\hat{L}(h)$.

Empirical Risk Minimization (ERM)

- Choose a prediction function \hat{h} which minimizes the loss on the training data within some **restricted** class of functions \mathcal{H} .

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i)).$$

- The average loss on the training data is called **empirical risk** $\hat{L}(h)$.
- \mathcal{H} can be: linear functions, polynomials, trees of a given depth, rules, linear combinations of trees, etc.⁴



⁴ T. Hastie, R. Tibshirani, and J. H. Friedman. *Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2003

Overfitting

Can you think of a “learning” algorithm which always have zero training error and poor test error?

Overfitting

Can you think of a “learning” algorithm which always have zero training error and poor test error?

“**Memorization**”: given a training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, choose:

$$h(\mathbf{x}) = \begin{cases} y_i & \text{if } \mathbf{x} = \mathbf{x}_i \text{ for some } i = 1, \dots, n, \\ 0 & \text{otherwise.} \end{cases}$$

Overfitting

Can you think of a “learning” algorithm which always have zero training error and poor test error?

“**Memorization**”: given a training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, choose:

$$h(\mathbf{x}) = \begin{cases} y_i & \text{if } \mathbf{x} = \mathbf{x}_i \text{ for some } i = 1, \dots, n, \\ 0 & \text{otherwise.} \end{cases}$$

It is **very easy** to drop training error down to 0 if our function class is flexible enough.

Overfitting

Can you think of a “learning” algorithm which always have zero training error and poor test error?

“**Memorization**”: given a training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, choose:

$$h(\mathbf{x}) = \begin{cases} y_i & \text{if } \mathbf{x} = \mathbf{x}_i \text{ for some } i = 1, \dots, n, \\ 0 & \text{otherwise.} \end{cases}$$

It is **very easy** to drop training error down to 0 if our function class is flexible enough.

This phenomenon is called **overfitting** to the data.

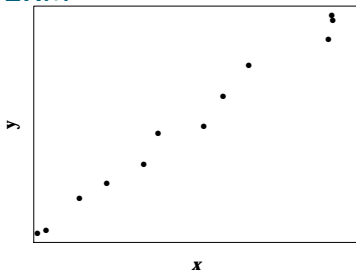
Overfitting – producing h with large **generalization error**

Overfitting with ERM

Simple artificial data set generator:

$$x \sim \text{uniform}(0, 1),$$

$$y = x + \epsilon, \quad \epsilon \sim N(0, 0.05).$$

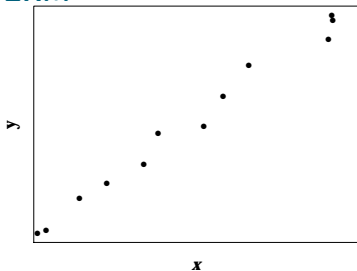


Overfitting with ERM

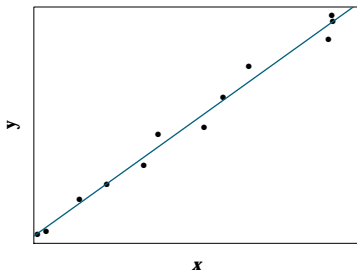
Simple artificial data set generator:

$$x \sim \text{uniform}(0, 1),$$

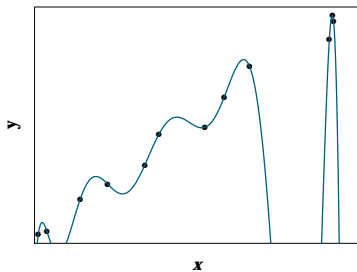
$$y = x + \epsilon, \quad \epsilon \sim N(0, 0.05).$$



Fitting linear function:



Fitting polynomial of degree n :



Overfitting

- A more complex/flexible function class is not always a good choice.
- Too complex class \implies overfitting.
- Too simple class \implies learning too constrained model.
- Can we control how much do we overfit?

Overfitting

- A more complex/flexible function class is not always a good choice.
- Too complex class \implies overfitting.
- Too simple class \implies learning too constrained model.
- Can we control how much do we overfit?

\implies **Vapnik-Chervonenkis theory** (only shown for 0/1 classification loss)

Generalization bounds

Theorem for finite classes (0/1 loss) — Occam's Razor⁵

Let the class of functions \mathcal{H} be finite. If the prediction function \hat{h} was trained on a sample of size n by minimizing the empirical risk within \mathcal{H} :

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{L}(h)$$

then with probability $1 - \delta$:

$$L(\hat{h}) - \hat{L}(\hat{h}) = O\left(\sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{n}}\right).$$

⁵ Alselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Occam's razor. *Inf. Process. Lett.*, 24(6):377–380, 1987

Generalization bounds

Theorem for finite classes (0/1 loss) — Occam's Razor⁵

Let the class of functions \mathcal{H} be finite. If the prediction function \hat{h} was trained on a sample of size n by minimizing the empirical risk within \mathcal{H} :

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{L}(h)$$

then with probability $1 - \delta$:

$$\underbrace{L(\hat{h}) - \hat{L}(\hat{h})}_{\text{generalization error}} = O\left(\sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{n}}\right).$$

⁵ Alselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K. Warmuth. Occam's razor. *Inf. Process. Lett.*, 24(6):377–380, 1987

Proof ⁶

A bin with red and green marbles.

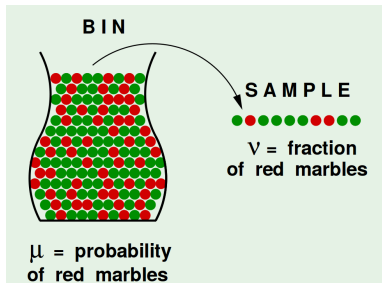
$$P(\text{picking a red marble}) = \mu,$$

$$P(\text{picking a green marble}) = 1 - \mu.$$

The value μ is unknown.

We pick n marbles **independently**.

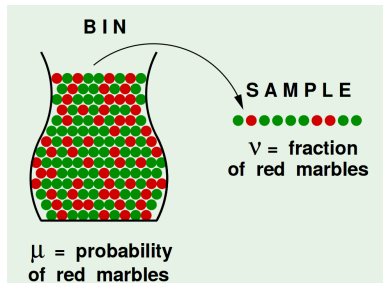
Fraction of red marbles in sample = ν .



⁶ Based on: Y. Abu-Mustafa, *Learning from data*, lecture notes at: <https://work.caltech.edu/lectures.html>

Proof

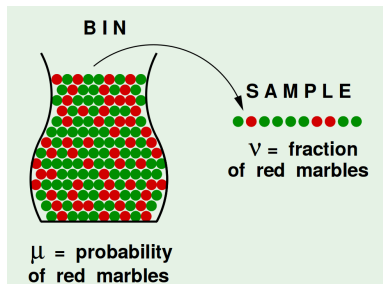
Does ν say anything about μ ?



Proof

Does ν say anything about μ ?

\implies **No:** ν Sample can be mostly green while bin is mostly red.

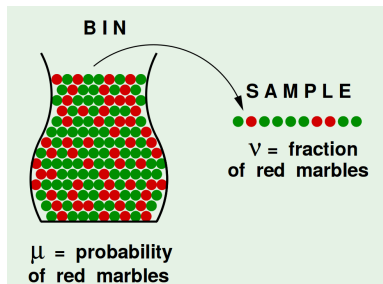


Proof

Does ν say anything about μ ?

⇒ **No:** ν Sample can be mostly **green** while bin is mostly **red**.

⇒ **Yes:** Sample frequency ν is **likely** to be close to bin frequency μ .



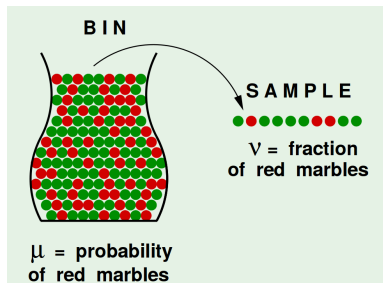
Proof

Does ν say anything about μ ?

Hoeffding's inequality:

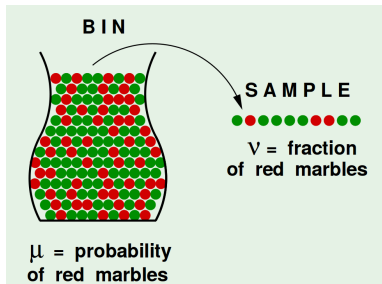
$$P(|\mu - \nu| > \epsilon) \leq 2e^{-2\epsilon^2 n}$$

Probability of large deviation of sample frequency from true probability decreases **exponentially** fast with n .



Proof

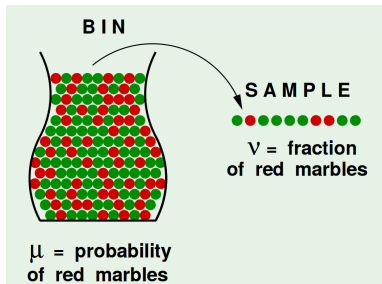
Given a **fixed** classifier h :



Proof

Given a **fixed** classifier h :

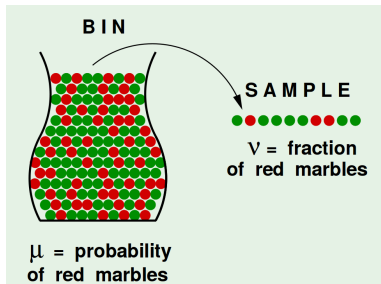
- Each example (x, y) is a marble



Proof

Given a **fixed** classifier h :

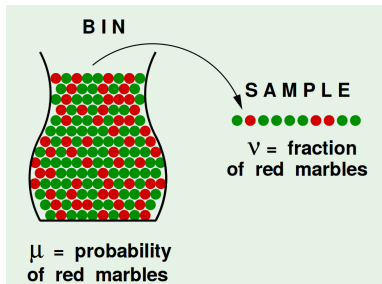
- Each example (\mathbf{x}, y) is a marble
- **red** marble = misclassification
($h(\mathbf{x}) \neq y$)
- **green** marble = correct classification
($h(\mathbf{x}) = y$)



Proof

Given a **fixed** classifier h :

- Each example (\mathbf{x}, y) is a marble
- **red** marble = misclassification
($h(\mathbf{x}) \neq y$)
- **green** marble = correct classification
($h(\mathbf{x}) = y$)

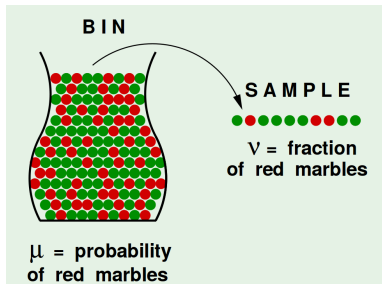


$$\mu = P(\text{picking a red marble}) = P(h \text{ misclassifies } (\mathbf{x}, y)) = L(h).$$

Proof

Given a **fixed** classifier h :

- Each example (\mathbf{x}, y) is a marble
- **red** marble = misclassification
($h(\mathbf{x}) \neq y$)
- **green** marble = correct classification
($h(\mathbf{x}) = y$)



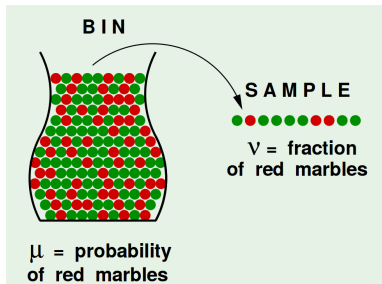
$$\mu = P(\text{picking a red marble}) = P(h \text{ misclassifies } (\mathbf{x}, y)) = L(h).$$

$$\nu = \text{fraction of red marbles in a sample} = \widehat{L}(h).$$

Proof

Given a **fixed** classifier h :

- Each example (\mathbf{x}, y) is a marble
- **red** marble = misclassification
($h(\mathbf{x}) \neq y$)
- **green** marble = correct classification
($h(\mathbf{x}) = y$)



$$\mu = P(\text{picking a red marble}) = P(h \text{ misclassifies } (\mathbf{x}, y)) = L(h).$$

$$\nu = \text{fraction of red marbles in a sample} = \widehat{L}(h).$$

Hoeffding's inequality: for any **fixed** classifier h ,

$$P(|\widehat{L}(h) - L(h)| > \epsilon) \leq 2e^{-2\epsilon^2 n}$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Given $\mathcal{H} = \{h_1, \dots, h_K\}$ is finite,

$$P(\exists h \in \mathcal{H}: |\widehat{L}(h) - L(h)| > \epsilon)$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Given $\mathcal{H} = \{h_1, \dots, h_K\}$ is finite,

$$\begin{aligned} P(\exists h \in \mathcal{H}: |\widehat{L}(h) - L(h)| > \epsilon) \\ = P\left(|\widehat{L}(h_1) - L(h_1)| > \epsilon \text{ or } \dots \text{ or } |\widehat{L}(h_K) - L(h_K)| > \epsilon\right) \end{aligned}$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Given $\mathcal{H} = \{h_1, \dots, h_K\}$ is finite,

$$\begin{aligned} &P(\exists h \in \mathcal{H}: |\widehat{L}(h) - L(h)| > \epsilon) \\ &= P\left(|\widehat{L}(h_1) - L(h_1)| > \epsilon \text{ or } \dots \text{ or } |\widehat{L}(h_K) - L(h_K)| > \epsilon\right) \\ \text{(union bnd)} \quad &= P(|\widehat{L}(h_1) - L(h_1)| > \epsilon) + \dots + P(|\widehat{L}(h_K) - L(h_K)| > \epsilon) \end{aligned}$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Given $\mathcal{H} = \{h_1, \dots, h_K\}$ is finite,

$$P(\exists h \in \mathcal{H}: |\widehat{L}(h) - L(h)| > \epsilon)$$

$$= P\left(|\widehat{L}(h_1) - L(h_1)| > \epsilon \text{ or } \dots \text{ or } |\widehat{L}(h_K) - L(h_K)| > \epsilon\right)$$

$$\text{(union bnd)} \quad = P(|\widehat{L}(h_1) - L(h_1)| > \epsilon) + \dots + P(|\widehat{L}(h_K) - L(h_K)| > \epsilon)$$

$$\text{(Hoeffding)} \quad \leq K \times 2e^{-2\epsilon^2 n}$$

Proof

How to handle the fact that h is **not fixed**, but **learns on a sample**?

Union bound

$$P(A_1 \text{ or } A_2 \text{ or } \dots \text{ or } A_K) \leq P(A_1) + P(A_2) + \dots + P(A_K).$$

Given $\mathcal{H} = \{h_1, \dots, h_K\}$ is finite,

$$P(\exists h \in \mathcal{H}: |\widehat{L}(h) - L(h)| > \epsilon)$$

$$= P\left(|\widehat{L}(h_1) - L(h_1)| > \epsilon \text{ or } \dots \text{ or } |\widehat{L}(h_K) - L(h_K)| > \epsilon\right)$$

$$\text{(union bnd)} \quad = P(|\widehat{L}(h_1) - L(h_1)| > \epsilon) + \dots + P(|\widehat{L}(h_K) - L(h_K)| > \epsilon)$$

$$\text{(Hoeffding)} \quad \leq K \times 2e^{-2\epsilon^2 n}$$

Probability that the training error of **any** classifier from \mathcal{H} deviates by more than ϵ from its test error is at most $2Ke^{-2\epsilon^2 n}$, where $K = |\mathcal{H}|$.

But:

$$\begin{aligned} P(\text{train. err. of } \hat{h} \text{ deviates from test err.}) \\ \leq P(\text{train. err. of any } h \text{ deviates from test. err.}) \end{aligned}$$

But:

$$\begin{aligned} P(\text{train. err. of } \hat{h} \text{ deviates from test err.}) \\ \leq P(\text{train. err. of any } h \text{ deviates from test. err.}) \end{aligned}$$

$$\begin{aligned} P(|\hat{L}(\hat{h}) - L(\hat{h})| > \epsilon) &\leq P(\exists h \in \mathcal{H}: |\hat{L}(h) - L(h)| > \epsilon) \\ &\leq \underbrace{2Ke^{-2\epsilon^2 n}}_{=\delta}. \end{aligned}$$

But:

$$\begin{aligned} P(\text{train. err. of } \hat{h} \text{ deviates from test err.}) \\ \leq P(\text{train. err. of any } h \text{ deviates from test. err.}) \end{aligned}$$

$$\begin{aligned} P(|\hat{L}(\hat{h}) - L(\hat{h})| > \epsilon) &\leq P(\exists h \in \mathcal{H}: |\hat{L}(h) - L(h)| > \epsilon) \\ &\leq \underbrace{2Ke^{-2\epsilon^2 n}}_{=\delta}. \end{aligned}$$

$$\delta = 2Ke^{-2\epsilon^2 n} \quad \iff \quad \epsilon = \sqrt{\frac{\log(K) + \log(2/\delta)}{2n}}$$

But:

$$\begin{aligned} P(\text{train. err. of } \hat{h} \text{ deviates from test err.}) \\ \leq P(\text{train. err. of any } h \text{ deviates from test. err.}) \end{aligned}$$

$$\begin{aligned} P(|\hat{L}(\hat{h}) - L(\hat{h})| > \epsilon) &\leq P(\exists h \in \mathcal{H}: |\hat{L}(h) - L(h)| > \epsilon) \\ &\leq \underbrace{2Ke^{-2\epsilon^2 n}}_{=\delta}. \end{aligned}$$

$$\delta = 2Ke^{-2\epsilon^2 n} \iff \epsilon = \sqrt{\frac{\log(K) + \log(2/\delta)}{2n}}$$

Conclusion: With probability at most $1 - \delta$:

$$L(\hat{h}) - \hat{L}(\hat{h}) \leq \epsilon = \sqrt{\frac{\log(|\mathcal{H}|) + \log(2/\delta)}{2n}}.$$

Generalization bounds

$$L(\hat{h}) = \hat{L}(\hat{h}) + O\left(\sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{n}}\right).$$

- Overhead grows very **slowly** (logarithmically) with size of $|\mathcal{H}|$.
- Overhead decreases as $O(1/\sqrt{n})$ with the sample size.
- Impractical: does not work for infinite classes of functions.

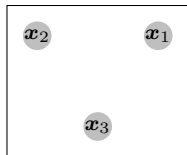
Vapnik-Chervonenkis (VC) dimension

A function class \mathcal{H} can **shatter** a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ if for any set of 0/1 labels $\{y_1, \dots, y_n\}$ there exists $h \in \mathcal{H}$ which gets training error 0.

Vapnik-Chervonenkis (VC) dimension

A function class \mathcal{H} can **shatter** a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ if for any set of 0/1 labels $\{y_1, \dots, y_n\}$ there exists $h \in \mathcal{H}$ which gets training error 0.

Example

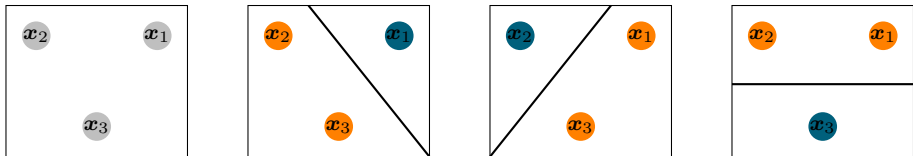


Class of linear functions in $d = 2$ can shatter any three points in general positions.

Vapnik-Chervonenkis (VC) dimension

A function class \mathcal{H} can **shatter** a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ if for any set of 0/1 labels $\{y_1, \dots, y_n\}$ there exists $h \in \mathcal{H}$ which gets training error 0.

Example



Class of linear functions in $d = 2$ can shatter any three points in general positions.

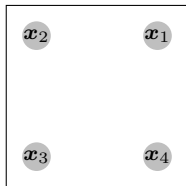
Vapnik-Chervonenkis (VC) dimension

The largest n , such that there exists a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ which can be shattered by \mathcal{H} is called **VC dimension** of \mathcal{H} .

Vapnik-Chervonenkis (VC) dimension

The largest n , such that there exists a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ which can be shattered by \mathcal{H} is called **VC dimension** of \mathcal{H} .

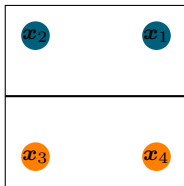
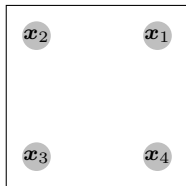
Example Class of linear functions in $d = 3$ has VC dimension 3.



Vapnik-Chervonenkis (VC) dimension

The largest n , such that there exists a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ which can be shattered by \mathcal{H} is called **VC dimension** of \mathcal{H} .

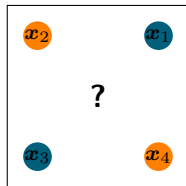
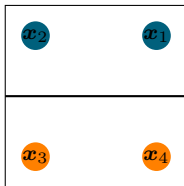
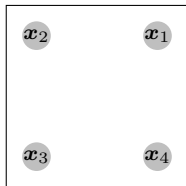
Example Class of linear functions in $d = 3$ has VC dimension 3.



Vapnik-Chervonenkis (VC) dimension

The largest n , such that there exists a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ which can be shattered by \mathcal{H} is called **VC dimension** of \mathcal{H} .

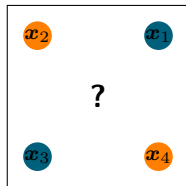
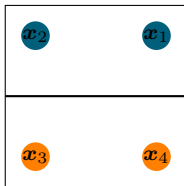
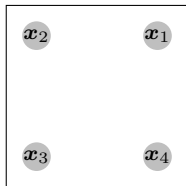
Example Class of linear functions in $d = 3$ has VC dimension 3.



Vapnik-Chervonenkis (VC) dimension

The largest n , such that there exists a set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ which can be shattered by \mathcal{H} is called **VC dimension** of \mathcal{H} .

Example Class of linear functions in $d = 3$ has VC dimension 3.



Generally, class of linear functions has VC dimension $d + 1$.

Vapnik-Chervonenkis (VC) dimension

What is VC dimension of:

- Circles centered at origin?
- Rectangles on a plane?
- Sine functions in one dimension?

Vapnik-Chervonenkis (VC) dimension

What is VC dimension of:

- Circles centered at origin? **VC-dim** = 2
- Rectangles on a plane?
- Sine functions in one dimension?

Vapnik-Chervonenkis (VC) dimension

What is VC dimension of:

- Circles centered at origin? **VC-dim** = 2
- Rectangles on a plane? **VC-dim** = 4
- Sine functions in one dimension?

Vapnik-Chervonenkis (VC) dimension

What is VC dimension of:

- Circles centered at origin? **VC-dim** = 2
- Rectangles on a plane? **VC-dim** = 4
- Sine functions in one dimension? **VC-dim** = ∞

Vapnik-Chervonenkis (VC) dimension

What is VC dimension of:

- Circles centered at origin? **VC-dim** = 2
- Rectangles on a plane? **VC-dim** = 4
- Sine functions in one dimension? **VC-dim** = ∞

In general, **VC-dim** \neq **num. of parameters!**

VC dimension measures the flexibility of a class of functions for classification purposes.

Generalization bounds

Theorem for VC classes (0/1 loss)⁷

Let the class of functions \mathcal{H} has VC dimension d_{VC} . If the prediction function \hat{h} was trained on a sample of size n by minimizing the empirical risk within \mathcal{H} :

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{L}(h)$$

then with probability $1 - \delta$:

$$\underbrace{L(\hat{h}) - \hat{L}(\hat{h})}_{\text{generalization error}} = O\left(\sqrt{\frac{d_{\text{VC}} + \log(1/\delta)}{n}}\right).$$



⁷ V. N. Vapnik and A. Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971

Generalization bounds – summary

- A theoretical framework behind machine learning which helps us to understand why and how learning algorithms work.
- Quantification of overfitting.
- The bounds should only be used qualitatively, not quantitatively! (they account for the worst case and are loose in general).
- Much more than what was presented here:
 - ▶ tighter bounds in special cases,
 - ▶ data dependent bounds (margin bounds, PAC-Bayes bounds),
 - ▶ other loss functions (regression, convex surrogates),
 - ▶ etc.

Estimation and approximation error

Estimation and approximation error

- If we learn h within class \mathcal{H} , we can only hope for getting h close to the **best prediction function within \mathcal{H}** :

$$h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} L(h)$$

Estimation and approximation error

- If we learn h within class \mathcal{H} , we can only hope for getting h close to the **best prediction function within \mathcal{H}** :

$$h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} L(h)$$

compare with $h^* = \arg \min_h L(h)$

Estimation and approximation error

- If we learn h within class \mathcal{H} , we can only hope for getting h close to the **best prediction function within \mathcal{H}** :

$$h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} L(h)$$

- Decomposition

compare with $h^* = \arg \min_h L(h)$

$$\begin{aligned} \text{Reg}(h) &= L(h) - L(h^*) \\ &= \underbrace{L(h) - L(h_{\mathcal{H}}^*)}_{\text{estimation error}} + \underbrace{L(h_{\mathcal{H}}^*) - L(h^*)}_{\text{approximation error}} \end{aligned}$$

Estimation and approximation error

- If we learn h within class \mathcal{H} , we can only hope for getting h close to the **best prediction function within \mathcal{H}** :

$$h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} L(h)$$

- Decomposition

compare with $h^* = \arg \min_h L(h)$

$$\begin{aligned} \text{Reg}(h) &= L(h) - L(h^*) \\ &= \underbrace{L(h) - L(h_{\mathcal{H}}^*)}_{\text{estimation error}} + \underbrace{L(h_{\mathcal{H}}^*) - L(h^*)}_{\text{approximation error}} \end{aligned}$$

- Only the estimation error depends on the training data.
- The approximation error depends on \mathcal{H} and $P(\mathbf{x}, y)$.

Estimation and approximation error

- If we learn h within class \mathcal{H} , we can only hope for getting h close to the **best prediction function within \mathcal{H}** :

$$h_{\mathcal{H}}^* = \arg \min_{h \in \mathcal{H}} L(h)$$

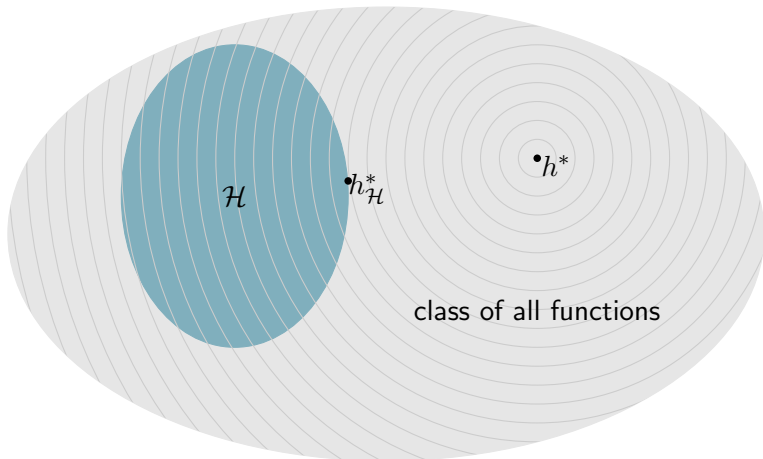
- Decomposition

compare with $h^* = \arg \min_h L(h)$

$$\begin{aligned} \text{Reg}(h) &= L(h) - L(h^*) \\ &= \underbrace{L(h) - L(h_{\mathcal{H}}^*)}_{\text{estimation error}} + \underbrace{L(h_{\mathcal{H}}^*) - L(h^*)}_{\text{approximation error}} \end{aligned}$$

- Only the estimation error depends on the training data.
- The approximation error depends on \mathcal{H} and $P(\mathbf{x}, y)$.
- A version of this decomposition: **bias-variance decomposition**.

Estimation and approximation error



Estimation and approximation error

