

# Introduction à l'apprentissage automatique - Polytech

## Empirical risk minimization

Caio Corro

Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique,  
91400, Orsay, France

# Theory

# Problem

## Train/test objective mismatch

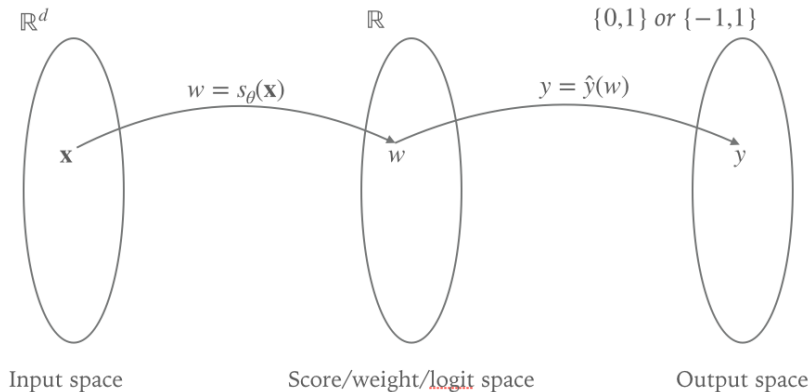
- ▶ At training time we minimize a loss function, e.g. the hinge loss
- ▶ At test time we evaluate the model using a different function, e.g. classification error / classification accuracy

There is a mismatch between the two objective!

## Goal

Can we prove that minimizing a given loss function also minimizes the test time objective?

# Assumptions



- ▶ We assume the class of scoring functions  $S$  is “rich enough” (*i.e.* set of all measurable mappings)
- ▶ We assume we have access to an infinite number of training datapoints

## 0-1 loss function

### Definition of the 0-1 loss function

Count the number of classification errors

⇒ exactly what we aim to minimize at test time

### Binary classification

- ▶ If  $Y = \{0, 1\}$ :  $\ell_{0-1}(y, w) = \mathbb{1}[(2y - 1)w < 0] = [\text{sign}(w) \neq 2y - 1]$
- ▶ If  $Y = \{-1, 1\}$ :  $\ell_{0-1}(y, w) = \mathbb{1}[y \times w < 0] = [\text{sign}(w) \neq y]$

### Multiclass classification

- ▶ If  $Y = E(k)$ :

$$\ell_{0-1}(\mathbf{y}, \mathbf{w}) = \begin{cases} 0 & \text{if } \mathbf{y} = \arg \max_{\mathbf{y}' \in E(k)} \langle \mathbf{w}, \mathbf{y}' \rangle, \\ 1 & \text{otherwise.} \end{cases}$$

- ▶ If  $Y = \{1, \dots, k\}$ :

$$\ell_{0-1}(y, \mathbf{w}) = \mathbb{1}[y \neq \arg \max_{y' \in \{1, \dots, k\}} w_{y'}]$$

# Bayes risk

## Notations

- ▶  $\mathbf{x}$ : random variable representing inputs in  $\mathbb{R}^d$
- ▶  $\mathbf{y}$ : random variable representing outputs, problem dependent
- ▶  $p(\mathbf{x}, \mathbf{y})$ : data distribution, unknown in practice but we can still use it for theory
- ▶  $S$ : set of scoring functions

## WARNING

The input/output mapping in the data distribution is not necessarily deterministic, an input  $\mathbf{x} \in \mathbb{R}^d$  may be associated with several outputs with a non null probability.

# Bayes risk

## Notations

- ▶  $\mathbf{x}$ : random variable representing inputs in  $\mathbb{R}^d$
- ▶  $\mathbf{y}$ : random variable representing outputs, problem dependent
- ▶  $p(\mathbf{x}, \mathbf{y})$ : data distribution, unknown in practice but we can still use it for theory
- ▶  $S$ : set of scoring functions

## WARNING

The input/output mapping in the data distribution is not necessarily deterministic, an input  $\mathbf{x} \in \mathbb{R}^d$  may be associated with several outputs with a non null probability.

## Risk of a scoring function

The risk of a given scoring function  $s \in S$  is denoted:

$$r(s) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \ell_{0-1}(\mathbf{y}, s(\mathbf{x})) ]$$

or, in other words, it is the classification error probability for classifier based on  $s$ .

# Minimum Bayes risk

## Risk of a scoring function

The risk of a given scoring function  $s \in S$  is denoted:

$$r(s) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \ell_{0-1}(\mathbf{y}, s(\mathbf{x})) ]$$

or, in other words, it is the classification error probability for classifier based on  $s$ .

## Minimum Bayes risk

It seems that it is a good idea to aim for a model of minimum Bayes risk:

$$\begin{aligned} r^* &= \inf_{s \in S} r(s) = \inf_{s \in S} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \ell_{0-1}(\mathbf{y}, s(\mathbf{x})) ] \\ &= \inf_{s \in S} \mathbb{E}_{\mathbf{x}}[ \mathbb{E}_{\mathbf{y}|\mathbf{x}}[ \ell_{0-1}(\mathbf{y}, s(\mathbf{x})) ] ] \\ &= \mathbb{E}_{\mathbf{x}}[ 1 - \max_{\mathbf{y} \in Y} p(\mathbf{y} = \mathbf{y}|\mathbf{x}) ] \end{aligned}$$



# Surrogate losses

## Recall

In practice we cannot use the 0-1 loss:

- ▶ non-convex
- ▶ partial derivatives are null almost everywhere

It is known that minimizing the 0-1 loss is a NP-hard problem (Ben-David et al., 2003; Feldman et al., 2009)

## Main idea

- ▶ Replace the 0-1 loss with a surrogate loss
- ▶ hope (prove?) that minimizing the surrogate loss leads to a classifier of minimum Bayes risk

# Classification calibration

## Surrogate risk

Let  $\tilde{\ell}$  be a surrogate loss.

The surrogate risk  $\tilde{r}(s)$  of scoring function  $s \in S$  is defined as:

$$\tilde{r}(s) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \tilde{\ell}(\mathbf{y}, s(\mathbf{x})) ],$$

and the optimal surrogate risk  $\tilde{r}^*$  is defined as:

$$\tilde{r}^* = \inf_{s \in S} \tilde{r}(s) = \inf_{s \in S} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \tilde{\ell}(\mathbf{y}, s(\mathbf{x})) ]$$

# Classification calibration

## Surrogate risk

Let  $\tilde{\ell}$  be a surrogate loss.

The surrogate risk  $\tilde{r}(s)$  of scoring function  $s \in S$  is defined as:

$$\tilde{r}(s) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \tilde{\ell}(\mathbf{y}, s(\mathbf{x})) ],$$

and the optimal surrogate risk  $\tilde{r}^*$  is defined as:

$$\tilde{r}^* = \inf_{s \in S} \tilde{r}(s) = \inf_{s \in S} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[ \tilde{\ell}(\mathbf{y}, s(\mathbf{x})) ]$$

## Definition

A surrogate loss  $\tilde{\ell}$  is said to be **classification calibrated** if and only if:

$$s^* \in \arg \min_{s \in S} \tilde{r}(s) \implies r(s^*) = r^*,$$

*i.e.* minimizing the surrogate risk leads to a prediction model of optimal Bayes risk.  
This property is also called Bayes consistency and Fisher consistency.

# Pointwise analysis 1/2

## Assumption

We assume the class of scoring function  $S$  is “rich enough”  
(i.e. set of all measurable mappings)

$$\begin{aligned} r^* &= \inf_{s \in S} r(s) = \inf_{s \in S} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \ell_{0-1}(s(\mathbf{x}), \mathbf{y}) ] \\ &= \inf_{\forall \mathbf{x} \in X: \mathbf{w}^{(\mathbf{x})} \in \mathbb{R}^k} \mathbb{E}_{\mathbf{x}} \left[ \mathbb{E}_{\mathbf{y}|\mathbf{x}} [ \ell_{0-1}(\mathbf{y}, \mathbf{w}^{(\mathbf{x})}) ] \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[ \underbrace{\inf_{\mathbf{w}^{(\mathbf{x})} \in \mathbb{R}^k} \mathbb{E}_{\mathbf{y}|\mathbf{x}} [ \ell_{0-1}(\mathbf{y}, \mathbf{w}^{(\mathbf{x})}) ]}_{\text{We can focus on this part only}} \right] \end{aligned}$$

# Pointwise analysis 1/2

## Assumption

We assume the class of scoring function  $S$  is “rich enough”  
(i.e. set of all measurable mappings)

## Pointwise risk

Let  $\mathbf{x} \in X$  such that  $p(\mathbf{x} = \mathbf{x}) > 0$ . We redefine the concept of (optimal) surrogate/Bayes risk as follows:

$$r(\mathbf{w}) = \mathbb{E}_{\mathbf{y}|\mathbf{x}=\mathbf{x}}[ \ell_{0-1}(\mathbf{y}, \mathbf{w}) ]$$

$$r^* = \inf_{\mathbf{w} \in \mathbb{R}^k} r(\mathbf{w})$$

$$\tilde{r}(\mathbf{w}) = \mathbb{E}_{\mathbf{y}|\mathbf{x}=\mathbf{x}}[ \tilde{\ell}(\mathbf{y}, \mathbf{w}) ]$$

$$\tilde{r}^* = \inf_{\mathbf{w} \in \mathbb{R}^k} \tilde{r}(\mathbf{w})$$

where  $\mathbf{w} \in \mathbb{R}^k$  should be interpreted as the output of the scoring function, i.e.  $\mathbf{w} = s(\mathbf{x})$ .

# Classification calibration for binary classification

## Surrogate losses

Let  $Y = \{-1, 1\}$ .

- ▶ Perceptron loss:  $\tilde{\ell}_{\text{perceptron}}(y, w) = \max(0, -y \times w)$
- ▶ Hinge loss:  $\tilde{\ell}_{\text{hinge}}(y, w) = \max(0, 1 - y \times w)$
- ▶ Squared hinge loss:  $\tilde{\ell}_{\text{s. hinge}}(y, w) = \max(0, 1 - y \times w)^2$
- ▶ Exponential loss:  $\tilde{\ell}_{\text{exp}}(y, w) = \exp(-y \times w)$
- ▶ Negative log-likelihood (NLL):  $\tilde{\ell}_{\text{nll}}(y, w) = \log(1 + \exp(-y \times w))$
- ▶ Quadratic error (or squared error):  $\tilde{\ell}_{\text{quad.}}(y, w) = (y \times w - 1)^2$

## Classification calibration

All these surrogate losses are classification calibrated **except the perceptron loss**.

# Classification calibration for binary classification

## Surrogate risk

$$\begin{aligned}\tilde{r}(w) &= \mathbb{E}_{\mathbf{y}|\mathbf{x}=\mathbf{x}}[ \tilde{\ell}(\mathbf{y}, w) ] \\ &= p(\mathbf{y} = 1|\mathbf{x} = \mathbf{x}) \times \tilde{\ell}(1, w) \quad + \quad p(\mathbf{y} = -1|\mathbf{x} = \mathbf{x}) \times \tilde{\ell}(-1, w) \\ &= \mu \times \tilde{\ell}(1, w) \quad + \quad (1 - \mu) \times \tilde{\ell}(-1, w)\end{aligned}$$

where  $\mu = p(\mathbf{y} = 1|\mathbf{x} = \mathbf{x})$ .

The optimal surrogate risk is:

$$\tilde{r}^* = \inf_{w \in \mathbb{R}} \tilde{r}(w) = \inf_{w \in \mathbb{R}} \mu \times \tilde{\ell}(1, w) + (1 - \mu) \times \tilde{\ell}(-1, w)$$

# Classification calibration for binary classification

## Surrogate risk

$$\begin{aligned}\tilde{r}(w) &= \mathbb{E}_{\mathbf{y}|\mathbf{x}=\mathbf{x}}[ \tilde{\ell}(\mathbf{y}, w) ] \\ &= p(\mathbf{y} = 1|\mathbf{x} = \mathbf{x}) \times \tilde{\ell}(1, w) \quad + \quad p(\mathbf{y} = -1|\mathbf{x} = \mathbf{x}) \times \tilde{\ell}(-1, w) \\ &= \mu \times \tilde{\ell}(1, w) \quad + \quad (1 - \mu) \times \tilde{\ell}(-1, w)\end{aligned}$$

where  $\mu = p(\mathbf{y} = 1|\mathbf{x} = \mathbf{x})$ .

The optimal surrogate risk is:

$$\tilde{r}^* = \inf_{w \in \mathbb{R}} \tilde{r}(w) = \inf_{w \in \mathbb{R}} \mu \times \tilde{\ell}(1, w) + (1 - \mu) \times \tilde{\ell}(-1, w)$$

## Intuition

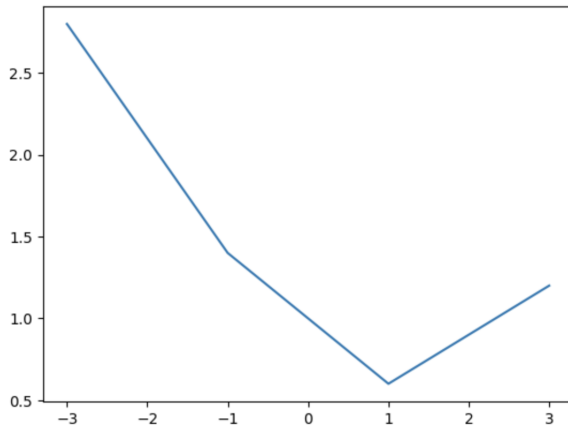
The surrogate loss  $\tilde{\ell}$  is classification calibrated if the minimizer  $w^*$  satisfies:

- ▶  $\mu > 0.5 \implies w^* \geq 0$  (the class 1 is the most probable class for input  $\mathbf{x}$ )
- ▶  $\mu < 0.5 \implies w^* < 0$  (the class 1 is the most probable class for input  $\mathbf{x}$ )
- ▶  $\mu = 0.5$ : this case is not important.



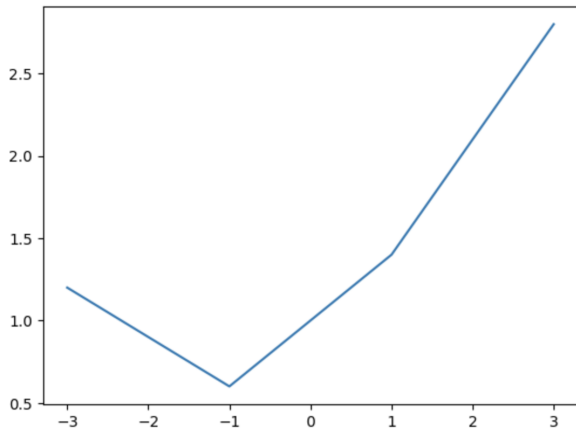
```
[4]: def hinge(y, w):  
      return max(0, 1 - y * w)  
  
      mu = 0.7  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * hinge(1, w) + (1 - mu) * hinge(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

```
[4]: [<matplotlib.lines.Line2D at 0x10f41f0a0>]
```



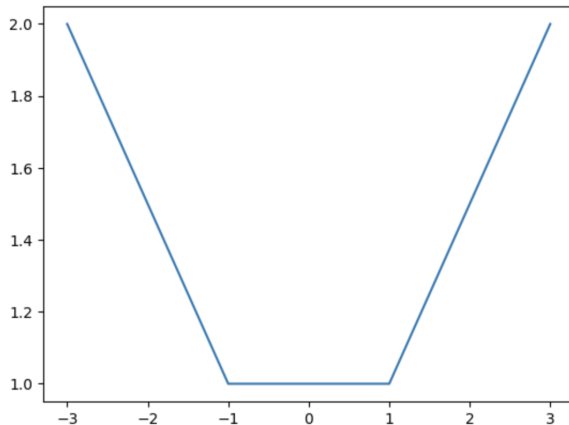
```
[5]: def hinge(y, w):  
      return max(0, 1 - y * w)  
  
      mu = 0.3  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * hinge(1, w) + (1 - mu) * hinge(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

[5]: [<matplotlib.lines.Line2D at 0x10f2313c0>]



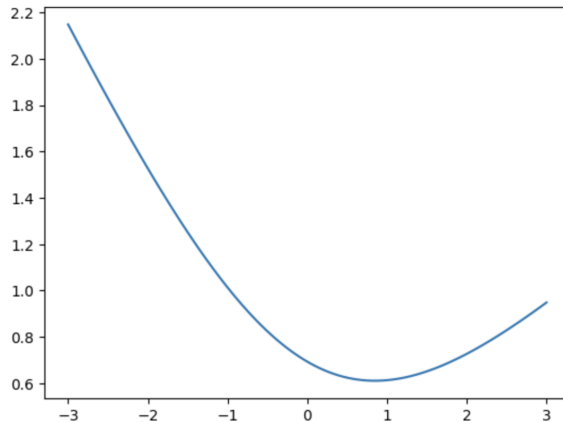
```
[6]: def hinge(y, w):  
      return max(0, 1 - y * w)  
  
      mu = 0.5  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * hinge(1, w) + (1 - mu) * hinge(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

```
[6]: [<matplotlib.lines.Line2D at 0x10ff3c820>]
```



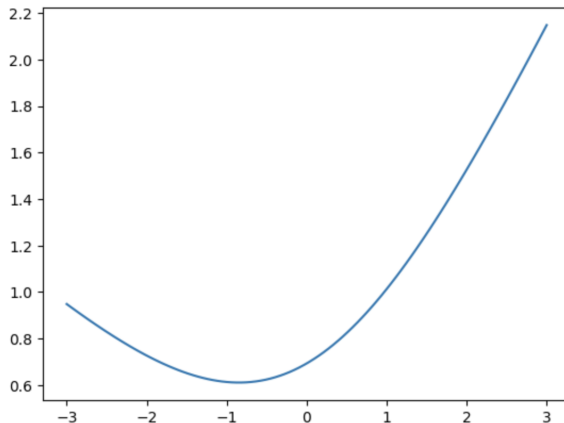
```
[7]: def nll(y, w):  
      return np.log(1 + np.exp(- y * w))  
  
      mu = 0.7  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * nll(1, w) + (1 - mu) * nll(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

```
[7]: [<matplotlib.lines.Line2D at 0x10ff9e290>]
```



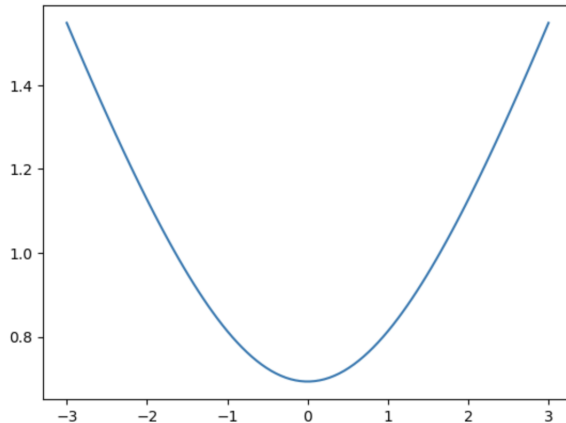
```
[8]: def nll(y, w):  
      return np.log(1 + np.exp(- y * w))  
  
      mu = 0.3  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * nll(1, w) + (1 - mu) * nll(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

```
[8]: [<matplotlib.lines.Line2D at 0x11002ded0>]
```



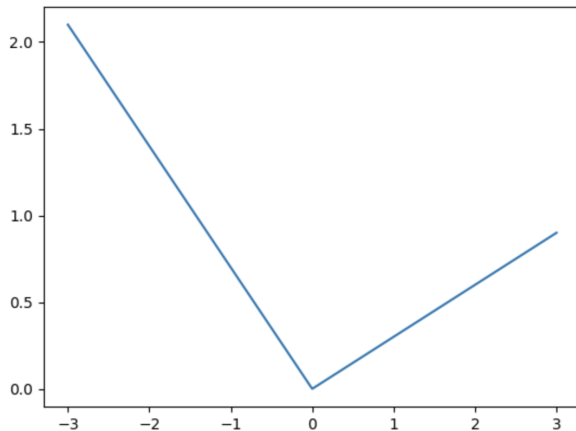
```
[9]: def nll(y, w):  
      return np.log(1 + np.exp(- y * w))  
  
      mu = 0.5  
      ws = np.linspace(-3, 3, 100)  
      ls = [mu * nll(1, w) + (1 - mu) * nll(-1, w) for w in ws]  
  
      plt.plot(ws, ls)
```

```
[9]: [<matplotlib.lines.Line2D at 0x1100b94b0>]
```



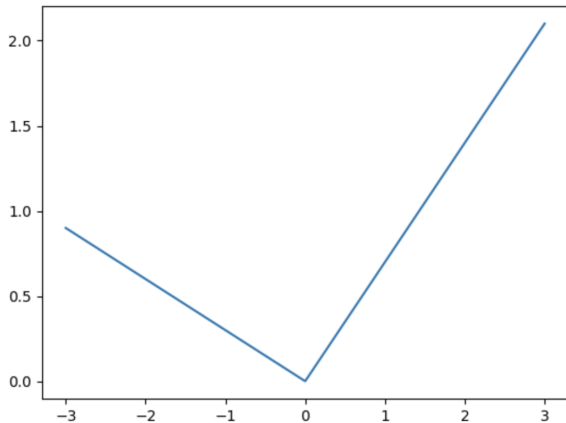
```
[14]: def perceptron(y, w):  
        return max(0, - y * w)  
  
        mu = 0.7  
        ws = np.linspace(-3, 3, 1000)  
        ls = [mu * perceptron(1, w) + (1 - mu) * perceptron(-1, w) for w in ws]  
  
        plt.plot(ws, ls)
```

```
[14]: [matplotlib.lines.Line2D at 0x110325a80>]
```



```
[15]: def perceptron(y, w):  
        return max(0, - y * w)  
  
        mu = 0.3  
        ws = np.linspace(-3, 3, 1000)  
        ls = [mu * perceptron(1, w) + (1 - mu) * perceptron(-1, w) for w in ws]  
  
        plt.plot(ws, ls)
```

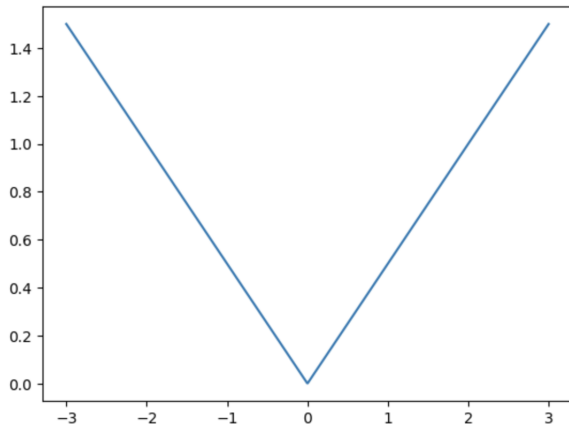
```
[15]: [matplotlib.lines.Line2D at 0x11039ad70<]
```





```
[16]: def perceptron(y, w):  
        return max(0, - y * w)  
  
        mu = 0.5  
        ws = np.linspace(-3, 3, 1000)  
        ls = [mu * perceptron(1, w) + (1 - mu) * perceptron(-1, w) for w in ws]  
  
        plt.plot(ws, ls)
```

```
[16]: [<matplotlib.lines.Line2D at 0x1104156f0>]
```



## Sufficient conditions in the binary case

Let  $\tilde{\ell} : \{-1, 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$  be a surrogate loss function that can be rewritten as:

$$\tilde{\ell}(y, w) = \phi(yw)$$

where  $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$  is function.

$\Rightarrow$  all previously presented binary loss function can be rewritten under this form.

### Theorem

If  $\phi$  is convex and differentiable at 0 with  $\phi'(0) < 0$ ,  
then  $\tilde{\ell}$  is classification calibrated

### Proof

See (Lin, 2004).

# Classification calibration for multiclass classification

## Surrogate losses

Let  $Y = E(k)$ .

- ▶ Hinge loss:  $\tilde{\ell}_{\text{hinge}}(\mathbf{y}, \mathbf{w}) = \max \left( 0, 1 - \langle \mathbf{y}, \mathbf{w} \rangle + \max_{\mathbf{y}' \in E(k) \setminus \{\mathbf{y}\}} \langle \mathbf{y}', \mathbf{w} \rangle \right)$
- ▶ NLL:  $\tilde{\ell}_{\text{nll}}(\mathbf{y}, \mathbf{w}) = -\langle \mathbf{y}, \mathbf{w} \rangle + \log \sum_i \exp(w_i)$

## Properties without proof

In the multiclass classification case:

- ▶ The hinge loss **is not** classification calibrated, see (Liu, 2007)
- ▶ The NLL loss **is** classification calibrated, see exercises

# Strictly proper losses 1/2

## Probabilistic prediction model

- ▶ We saw that we can learn models that predict a probability distribution over outputs, e.g.  $p_s(\mathbf{y}|\mathbf{x})$ , where the  $s$  emphasize that this is the learned model distribution, parameterized by the scoring function  $s$ .
- ▶ Classification calibration means the the most probable output in the data distribution will also be the most probable output in the model distribution
- ▶ We may want a stronger property: that the two distributions are equal

# Strictly proper losses 1/2

## Probabilistic prediction model

- ▶ We saw that we can learn models that predict a probability distribution over outputs, e.g.  $p_s(\mathbf{y}|\mathbf{x})$ , where the  $s$  emphasize that this is the learned model distribution, parameterized by the scoring function  $s$ .
- ▶ Classification calibration means the the most probable output in the data distribution will also be the most probable output in the model distribution
- ▶ We may want a stronger property: that the two distributions are equal

## Definition

Let  $p(\mathbf{y}, \mathbf{x})$  be the data distribution and  $p_s(\mathbf{y}|\mathbf{x})$  the model distribution. In the pointwise setting, a surrogate loss  $\tilde{\ell}$  is **classification calibrated** if and only if the scoring function  $s^*$  that minimizes the surrogate risk leads to a model distribution equal to the data distribution:

$$\forall \mathbf{y} \in Y : \quad p_{s^*}(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}) = p(\mathbf{y} = \mathbf{y} | \mathbf{x} = \mathbf{x}).$$

# Strictly proper losses 1/2

## Remarks

- ▶ Strict properness implies classification calibration
- ▶ The support of the data distribution must be “representable” by the model distribution
- ▶ The NLL loss is strictly proper for models whose probability parameters are computed using the sigmoid/softmax function **if the conditional data distribution has full support.**

# Risk minimization decomposition

## Practical issues

- ▶ We only have access to a finite training dataset
- ▶ The set of function  $S$  is not the set of all measurable mapping
- ▶ The learning algorithm may not find the optimal classifier  $s \in S$ , *i.e.* the following problem may be solved approximately

$$s^* \in \arg \min_{s \in S} \frac{1}{|D|} \sum_{(\mathbf{x}, \mathbf{y}) \in D} \ell(\mathbf{y}, s(\mathbf{x})),$$

or in other words, in practice the computed  $s^*$  is not a minimizer.

# Risk minimization decomposition

## Practical issues

- ▶ We only have access to a finite training dataset
- ▶ The set of function  $S$  is not the set of all measurable mapping
- ▶ The learning algorithm may not find the optimal classifier  $s \in S$ , *i.e.* the following problem may be solved approximately

$$s^* \in \arg \min_{s \in S} \frac{1}{|D|} \sum_{(\mathbf{x}, \mathbf{y}) \in D} \ell(\mathbf{y}, s(\mathbf{x})),$$

or in other words, in practice the computed  $s^*$  is not a minimizer.

## Risk decomposition

Excess risk:

$$r(s^*) - r^* \geq 0$$

Excess risk decomposition:

$$r(s^*) - r^* = \underbrace{r(s^*) - \inf_{s \in S} r(s)}_{\text{Estimation error}} + \underbrace{\inf_{s \in S} r(s) - r^*}_{\text{Approximation error}}$$