

# Sample Compressed SVM

*(or A sample compressed PAC-Bayes approach to kernel methods)*

Pascal Germain

Joint work with François Laviolette, Alexandre Lacasse,  
Alexandre Lacoste, Mario Marchand and Sara Shanian

GRAAL  
(Université Laval, Québec city)

February 18, 2010

In this lecture, we will :

- Review quickly the **Sample-Compress theory**
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers** (the Sc-SVM)
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- Present some experimental **results**
- and Conclude...

# The Classification problem

We consider a training set  $S$  of  $m$  examples

$$S \stackrel{\text{def}}{=} (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m)$$

where each  $\mathbf{z}_i$  is a input-output pair:

$$\mathbf{z}_i \stackrel{\text{def}}{=} (\mathbf{x}_i, y_i)$$

$$\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^n \quad (\text{Real attributes})$$

$$y_i \in \mathcal{Y} = \{-1, +1\} \quad (\text{Binary classif.})$$

Each example  $\mathbf{z}_i$  is drawn *IID* according to an unknown probability distribution  $D$  on  $\mathcal{X} \times \mathcal{Y}$ . Hence :

$$S \sim D^m$$

# Elements of the Sample Compression theory

A **sc-classifier**  $h_i^\mu$  is a data-dependent classifier described by two variables:

- A **compression-set**  $S_i$  containing a subset of the training sequence  $S$  describing the classifier
  - $\mathbf{i} \stackrel{\text{def}}{=} \langle i_1, i_2, \dots, i_m \rangle$  with  $1 \leq i_1 < i_2 < \dots < i_{|\mathbf{i}|} \leq m$
- A **message string**  $\mu$  containing the additional information needed to construct the classifier.
  - $\mu$  is chosen among  $\mathcal{M}_i$ , a predefined set of all messages that can be supplied with  $S_i$ .

Given  $S_i$  and  $\mu$ , a **reconstruction function**  $\mathcal{R}$  outputs a classifier :

$$h_i^\mu \stackrel{\text{def}}{=} \mathcal{R}(S_i, \mu).$$

# Risk of a sc-classifier

The **risk** (or generalization error) of a classifier  $h$  is defined as

$$R_D(h) \stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y) \sim D} I(h(\mathbf{x}) \neq y) = \Pr_{(\mathbf{x}, y) \sim D} (h(\mathbf{x}) \neq y)$$

where  $I(a) = 1$  if predicate  $a$  is true and 0 otherwise.

The **empirical risk** of a sc-classifier  $h_i^\mu$  on the training set  $S$  is defined by

$$R_S(h_i^\mu) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j=1}^m R_{\langle (\mathbf{x}_j, y_j) \rangle} (h_i^\mu),$$

where

$$R_{\langle (\mathbf{x}_j, y_j) \rangle} (h_i^\mu) \stackrel{\text{def}}{=} \begin{cases} I(h_i^\mu(\mathbf{x}_j) \neq y_j) & \text{if } j \notin \mathbf{i} \\ 0 & \text{otherwise.} \end{cases}$$

Thus,  $mR_S(h_i^\mu) \sim \text{Bin}(m - \|\mathbf{i}\|, R_D(h_i^\mu))$ .

# Examples of compression sets and reconstruction functions

## Support Vector Machine

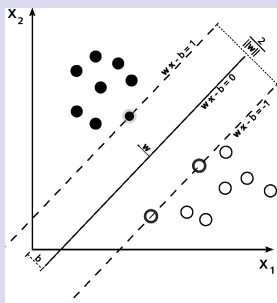


Image: Wikipedia

We can reconstruct a SVM by using the **support-vectors** as the compression-set.

In this example :

- Training set size:  $|S| = 16$
- Compression-set size:  $|i| = 3$
- Message string:  $\mu = \emptyset$

## Perceptron

Once a perceptron is trained (until convergence) on  $S$ , we only need the  $|i|$  **examples implied in an update** to reconstruct the classifier.

- Again,  $\mu = \emptyset$ .

In this lecture, we will :

- Review quickly the Sample-Compress theory
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers** (the Sc-SVM)
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- Present some experimental **results**
- and Conclude...

- In an attempt to **upper-bound the risk**  $R_D(h)$ , usual sample-compression bounds **degrade with the size** of the compression-set  $S_i$  (as we expect).
- This gives a **bad interpretation of the generalization error** of the SVM, which can have a low risk even if there is a large number of support vectors.
- To overcome this issue, let's define the SVM as a **majority vote of sc-classifiers** of unitary compression-size.



# Redefining the SVM

We denote  $\mathcal{H}^S$  the set of all sc-classifiers. Each  $h_i^\mu \in \mathcal{H}^S$  is such as :

- The **compression-set** contains only zero or one training example :

$$S_i \in \{S_{\langle \rangle}, S_{\langle 1 \rangle}, S_{\langle 2 \rangle}, \dots, S_{\langle m \rangle}\}$$

- The **message string** is formed by a real number and a sign :

$$\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$$

We have  $\mathcal{M}_{\langle i \rangle} = [-1, 1] \times \{+, -\}$  and  $\mathcal{M}_{\langle \rangle} = \{\epsilon\} \times \{+, -\}$ .

We consider **pairs of boolean complement classifiers** such as :

$$h_i^{(\sigma, -)}(\mathbf{x}) = -h_i^{(\sigma, +)}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}, \sigma \in [-1, 1].$$

We also have:

$$h_{\langle \rangle}^{(\epsilon, +)}(\mathbf{x}) = +1 \quad \text{and} \quad h_{\langle \rangle}^{(\epsilon, -)}(\mathbf{x}) = -1 \quad \forall \mathbf{x} \in \mathcal{X}.$$

## sc-classifier $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|I|} \times \{+, -\}$

## Distribution $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(\mathbf{i})Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma,+)} - Q(h_i^{(\sigma,-)}) = w_i$$

Let  $Q$  be a **probability distribution** over  $\mathcal{H}^S$ . We denote

- $Q_{\mathcal{I}}$ , the probability that a compression-set  $S_i$  is chosen by  $Q$ :

$$Q_{\mathcal{I}}(\mathbf{i}) \stackrel{\text{def}}{=} \int_{\mu \in \mathcal{M}_i} Q(h_i^\mu) d\mu$$

- $Q_{S_i}$ , the probability of choosing message  $\mu$  given  $S_i$ :

$$Q_{S_i}(\mu) \stackrel{\text{def}}{=} Q(h_i^\mu | S_i)$$

- Therefore,  $Q(h_i^\mu) = Q_{\mathcal{I}}(\mathbf{i})Q_{S_i}(\mu)$ .

The **output** of the majority vote classifier (*bayes classifier*) is given by :

$$B_Q(\mathbf{x}) \stackrel{\text{def}}{=} \text{sgn} \left[ \mathbf{E}_{h \sim Q} h(\mathbf{x}) \right]$$

## sc-classifier $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle 1 \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$

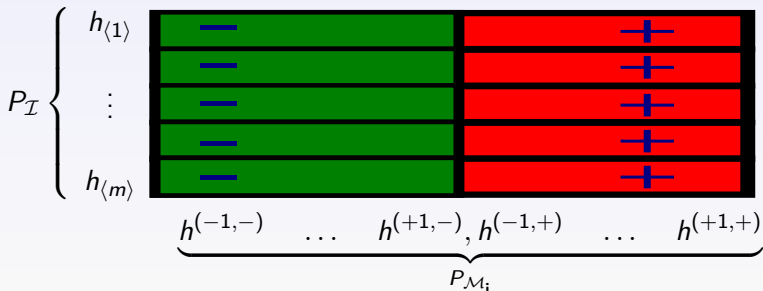
## Distribution $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(i) Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma, +)}) - Q(h_i^{(\sigma, -)}) = w_i$$

Before seeing the data, we define a **prior distribution** over the compression-sets and the message strings. This gives us indirectly a prior  $P$  over  $\mathcal{H}^S$  such as :

- $P_{\mathcal{I}}$  is an uniform distribution over all possible compression-sets ;
- For each compression-set  $S_i$ ,  $P_{S_i}$  is uniform over all messages.



## sc-classifier $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$

## Distribution $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(\mathbf{i})Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma,+)}) - Q(h_i^{(\sigma,-)}) = w_i$$

We say that a posterior  $Q$  is **aligned on** a prior  $P$  when for all  $\mathbf{i}$  and  $\sigma$ :

$$Q(h_i^{(\sigma,+)}) + Q(h_i^{(\sigma,-)}) = P(h_i^{(\sigma,+)}) + P(h_i^{(\sigma,-)})$$

Moreover, we say that a posterior  $Q$  is **strongly aligned** when for all  $\mathbf{i}$ , there is a  $w_i$  such that for all  $\sigma$ :

$$Q(h_i^{(\sigma,+)}) - Q(h_i^{(\sigma,-)}) = w_i$$

By restricting ourself to strongly aligned posterior, we obtain a posterior distribution totally defined by the  $w_i$ 's :

$$Q(h_i^{(\sigma,+)}) = \frac{1}{2} \left( P(h_i^{(\sigma,+)}) + P(h_i^{(\sigma,-)}) + w_i \right)$$

$$Q(h_i^{(\sigma,-)}) = \frac{1}{2} \left( P(h_i^{(\sigma,+)}) + P(h_i^{(\sigma,-)}) - w_i \right)$$

sc-classifier  $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle 1 \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$

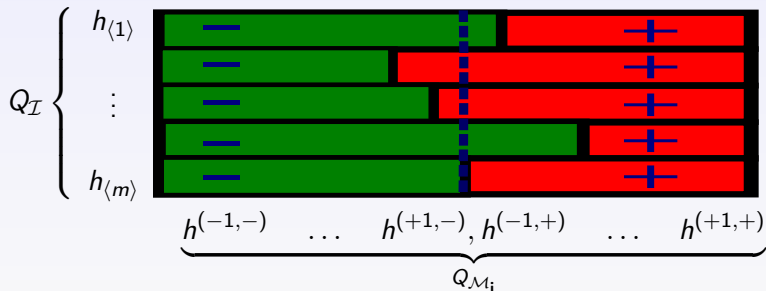
Distribution  $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(i) Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma,+)} - Q(h_i^{(\sigma,-)}) = w_i$$

$$Q(h_i^{(\sigma,+)} = \frac{1}{2} \left( P(h_i^{(\sigma,+)} + P(h_i^{(\sigma,-)}) + w_i \right)$$

$$Q(h_i^{(\sigma,-)} = \frac{1}{2} \left( P(h_i^{(\sigma,+)} + P(h_i^{(\sigma,-)}) - w_i \right)$$



sc-classifier  $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$

Distribution  $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(\mathbf{i}) Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma,+)}) - Q(h_i^{(\sigma,-)}) = w_i$$

There's almost **no loss of expressiveness** if we consider aligned posterior:

### Proposition

Let  $P$  be a prior,  $S$  a training sequence, and  $Q$  a distribution on  $\mathcal{H}^S$  for which there exists  $A > 0$  such that for all  $h_i^\mu$  :

$$Q(h_i^\mu) + Q(-h_i^\mu) \leq A(P(h_i^\mu) + P(-h_i^\mu)).$$

Then there exists a distribution  $Q'$  aligned on  $P$  and Bayes-equivalent to  $Q$

$$\text{i.e., } B_{Q'}(\mathbf{x}) = B_Q(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X}.$$

Remember that  $B_Q(\mathbf{x}) \stackrel{\text{def}}{=} \text{sgn} \left[ \mathbf{E}_{h \sim Q} h(\mathbf{x}) \right]$ .

From  $Q$  to  $Q'$ , the margins will vary, but the outcome of the majority vote will stay the same!

## sc-classifier $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle \cdot \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|i|} \times \{+, -\}$

## Distribution $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(\mathbf{i}) Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma,+)}) - Q(h_i^{(\sigma,-)}) = w_i$$

Consider any similarity function  $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow [-1, 1]$ .

We say that **reconstruction function**  $\mathcal{R}$  is associated to  $k$  when :

$$h_{\langle \cdot \rangle}^{(\epsilon, +)}(\mathbf{x}) \stackrel{\text{def}}{=} +1$$

$$h_{\langle i \rangle}^{(\sigma, +)}(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} +1 & \text{if } \sigma < k(\mathbf{x}_i, \mathbf{x}) \\ -1 & \text{otherwise} \end{cases}$$

$$h_i^{(\sigma, -)}(\mathbf{x}) \stackrel{\text{def}}{=} -h_i^{(\sigma, +)}(\mathbf{x}).$$

For an uniform prior, this definition allows us to **recover the value** of  $k$ :

$$\frac{1}{2} k(\mathbf{x}_i, \mathbf{x}) = \int_{\sigma \in \mathcal{M}^1(S_i)} h_i^{(\sigma, +)}(\mathbf{x}) \cdot P_{S_i}(\sigma, +) d\mu,$$

$$\text{as } \int_{-1}^{k(\mathbf{x}_i, \mathbf{x})} \frac{1/2}{(1-(-1))} d\mu - \int_{k(\mathbf{x}_i, \mathbf{x})}^1 \frac{1/2}{(1-(-1))} d\mu = \frac{1}{2} k(\mathbf{x}_i, \mathbf{x}).$$

## sc-classifier $h_i^\mu \in \mathcal{H}^S$

Comp-set:  $S_i \in \{S_{\langle \rangle}, S_{\langle 1 \rangle}, \dots, S_{\langle m \rangle}\}$

Message:  $\mu \in \mathcal{M}_i = [-1, 1]^{|I|} \times \{+, -\}$

## Distribution $Q$

$$Q(h_i^\mu) = Q_{\mathcal{I}}(i) Q_{S_i}(\mu)$$

$$Q(h_i^{(\sigma, +)}) - Q(h_i^{(\sigma, -)}) = w_i$$

We finally obtain that our strongly aligned posterior will be such that:

$$Q_{\mathcal{I}}(\langle \rangle) = Q_{\mathcal{I}}(\langle i \rangle) = \frac{1}{m+1},$$

$$w_{\langle i \rangle} \cdot k(\mathbf{x}_i, \mathbf{x}) = \int_{\mu \in \mathcal{M}_{\langle i \rangle}} h_{\langle i \rangle}^\mu(\mathbf{x}) \cdot Q_{\langle i \rangle}(\mu) d\mu,$$

$$\text{and } w_{\langle \rangle} \cdot 1 = \int_{\mu \in \mathcal{M}_{\langle \rangle}} h_{\langle \rangle}^\mu(\mathbf{x}) \cdot Q_{\langle \rangle}(\mu) d\mu.$$

Thus, the output of this majority vote  $B_Q(\mathbf{x}) = \text{sgn} \left[ \mathbf{E}_{h \sim Q} h(\mathbf{x}) \right]$  will be the same as  $f_{\text{SVM}}(\mathbf{x}) = \text{sgn} \left( \sum_{i=1}^m y_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b \right)$  when

$$w_{\langle i \rangle} = \frac{y_i \alpha_i}{Z(m+1)} \quad \text{and} \quad w_{\langle \rangle} = \frac{b}{Z(m+1)}. \quad \left( Z \stackrel{\text{def}}{=} \sum_{i=1}^m \alpha_i + |b| \right)$$



In this lecture, we will :

- Review quickly the Sample-Compress theory
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers** (the Sc-SVM)
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- Present some experimental **results**
- and Conclude...

# PAC-Bayes bounds for Sc-SVM

Usuals PAC-Bayes theorems allow us to bound the risk of a majority vote classifier using two key ingredients:

- The **Kullback-Leibler divergence**  $KL(Q||P)$  between prior distribution  $P$  and posterior distribution  $Q$
- The **empirical risk of the Gibbs classifier**  $G_Q$ , related to the majority vote  $B_Q$ 
  - Given any  $\mathbf{x}$ ,  $G_Q$  draws  $h$  according to  $Q$  and classifies  $\mathbf{x}$  according to  $h$ .
  - It follows that  $R_D(B_Q) \leq 2R_D(G_Q)$  .

In our setting, the Gibbs risk  $R_D(G_Q)$  will be likely near  $1/2$ , even if the Bayes risk is close to 0.

- Each sc-classifier  $h_i^\mu \in \mathcal{H}^S$  might be really weak.

We want to bound a **more relevant risk!**

Inspired by [Germain et al. *PAC-Bayes bounds for general loss functions* (2006)].

## Margin of the majority vote classifier

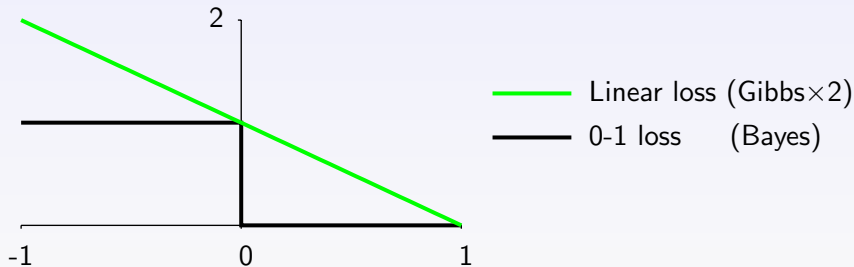
$$M_Q(\mathbf{x}, y) \stackrel{\text{def}}{=} \mathbf{E}_{h_i^\mu \sim Q} y h_i^\mu(\mathbf{x})$$

The margin is closely related to the Gibbs risk :

$$R_D(G_Q) = \frac{1}{2} - \frac{1}{2} \mathbf{E}_{(\mathbf{x}, y) \sim D} M_Q(\mathbf{x}, y).$$

For sc-classifiers, we define the **empirical margin** as:

$$\widehat{M}_Q(\mathbf{x}_j, y_j) \stackrel{\text{def}}{=} \mathbf{E}_{h_i^\mu \sim Q} \left[ y_j h_i^\mu(\mathbf{x}_j) \cdot I(j \notin \mathbf{i}) + 1 \cdot I(j \in \mathbf{i}) \right].$$



## Margin of the majority vote classifier

$$M_Q(\mathbf{x}, y) \stackrel{\text{def}}{=} \mathbf{E}_{h_i^\mu \sim Q} y h_i^\mu(\mathbf{x})$$

We consider **any non-negative loss**  $\zeta$  that can be expanded by a Taylor series around  $M_Q(\mathbf{x}, y) = 0$  and **upper-bound the zero-one loss**:

$$\zeta(\alpha) = \sum_{k=0}^{\text{deg } \zeta} a_k \alpha^k \quad \text{with } a_k \geq 0$$
$$\zeta(\alpha) \geq I(\alpha \leq 0) \quad \forall \alpha \in [-1, 1].$$

We obtain a risk value

$$\zeta_D^Q \stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y) \sim D} \sum_{k=0}^{\text{deg } \zeta} a_k (-M_Q(\mathbf{x}, y))^k$$

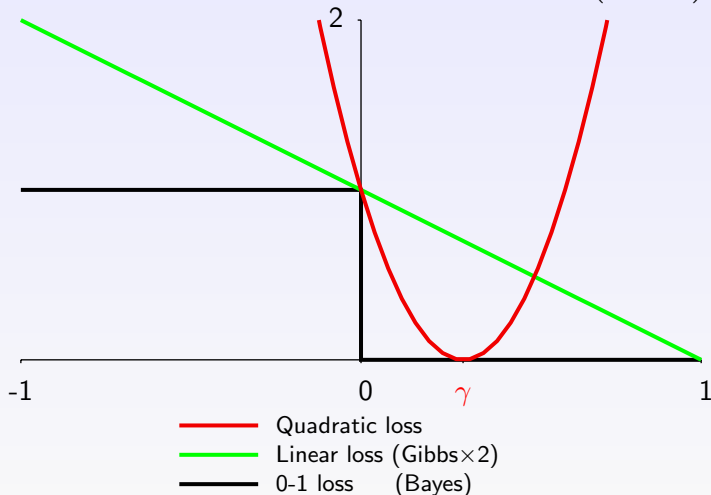
We express  $\zeta_D^Q$  as the risk of a Gibbs classifier described by a transformed posterior  $\bar{Q}$  on an augmented set of classifiers  $\bar{\mathcal{H}}^S$ :

$$\zeta_D^Q = \zeta(1) \cdot \mathbf{E}_{(\mathbf{x}, y) \sim D} \left[ \frac{1}{2} - \frac{1}{2} M_{\bar{Q}}(\mathbf{x}, y) \right]$$

## Margin of the majority vote classifier

$$M_Q(\mathbf{x}, y) \stackrel{\text{def}}{=} \mathbf{E}_{h_i^\mu \sim Q} y h_i^\mu(\mathbf{x})$$

We choose to use the **quadratic loss** function  $\zeta_\gamma(\alpha) = \left(1 - \frac{1}{\gamma}\alpha\right)^2$ .



# First PAC-Bayes theorem

The following PAC-Bayes theorem (next slide) is an adapted version of a **Catoni's theorem** where the influence of the empirical risk (vs  $KL(Q||P)$ ) is determined by an hyperparameter  $C$ .

- Generally less tight than the classic kl bound
- But useful to design a bound-minimization algorithm

In this adapted version, we consider:

- A general loss function  $\zeta$
- A set of (data-dependent) sc-classifiers of size  $\leq l$

## Theorem

For any  $D$ , any family  $(\mathcal{H}^S)_{S \in \mathcal{D}^m}$  of sets of sc-classifiers of size at most  $l$ , any prior  $P$  any  $\delta \in (0, 1]$ , any positive real number  $C_1$  and any margin loss function  $\zeta$  of degree  $< m/l$ , we have

$$\Pr_{S \sim D^m} \left( \forall Q \text{ on } \mathcal{H}^S: \zeta_S^Q \leq C' \cdot \left( \zeta_S^Q + \frac{\zeta'(1) \cdot \text{KL}(Q \| P) + \ln \frac{1}{\delta}}{\zeta(1) \cdot C_1 \cdot m} \right) \right) \geq 1 - \delta,$$

where  $\text{KL}(\cdot \| \cdot)$  is the Kullback-Leibler divergence, and

$$C' = \frac{C_1 \cdot \frac{m}{m - l \cdot \text{deg } \zeta}}{1 - e^{-C_1 \cdot \frac{m - l \cdot \text{deg } \zeta}{m}}}.$$

Finding  $Q$  that minimizes this bound is equivalent to finding  $Q$  that minimizes :

$$f(Q) \stackrel{\text{def}}{=} C \cdot \zeta_S^Q + \text{KL}(Q \| P)$$

## Second PAC-Bayes theorem

The next theorem is an adapted version of the **Langford and Seeger's theorem** where the influence of the empirical risk (vs  $\text{KL}(Q\|P)$ ) is given via the Kullback-Leibler divergence between two Bernoulli distributions of probability of success  $p$  and  $q$  :

$$\begin{aligned}\text{kl}(q\|p) &\stackrel{\text{def}}{=} q \ln \frac{q}{p} + (1 - q) \ln \frac{1 - q}{1 - p} \\ &= \text{kl}(1 - q\|1 - p)\end{aligned}$$

We specialize the theorem for the case of aligned posterior:

$$Q(h) + Q(-h) = P(h) + P(-h) \quad \forall h \in \mathcal{H}$$

...And the term  $\text{KL}(Q\|P)$  disappears from the theorem!



## Theorem

For any  $D$ , any family  $(\mathcal{H}^S)_{S \in \mathcal{D}^m}$  of sets of sc-classifiers of size at most  $l$ , any prior  $P$ , any  $\delta \in (0, 1]$ , any margin loss function  $\zeta$  of degree  $< m/l$ , we have

$$\Pr_{S \sim D^m} \left( \forall Q \in \mathcal{H}^S \text{ aligned on } P: \text{kl} \left( \frac{1}{\zeta(1)} \cdot \zeta_S^Q \parallel \frac{1}{\zeta(1)} \cdot \zeta_D^Q \right) \leq \frac{\ln \frac{m+1}{\delta}}{m-l \cdot \text{deg } \zeta} \right) \geq 1 - \delta$$

where  $\text{kl}(q \parallel p)$  is the KL-divergence between two Bernoulli distributions of respective success probabilities  $q$  and  $p$ .

Finding  $Q$  that minimizes this bound is equivalent to finding  $Q$  that minimizes :

$$f(Q) \stackrel{\text{def}}{=} \zeta_S^Q$$

We want to bound random variable  $\mathbf{E}_{h \sim P} e^{m \cdot \text{kl}(R_S(h) \| R(h))}$  in term of  $R(G_Q)$ .

### General theorem

Term  $\text{KL}(Q \| P)$  arises when transforming expectation over  $P$  into expectation over  $Q$ :

$$\begin{aligned}
 & \ln \left[ \mathbf{E}_{h \sim P} e^{m \cdot \text{kl}(R_S(h) \| R(h))} \right] \\
 &= \ln \left[ \mathbf{E}_{h \sim Q} \frac{P(h)}{Q(h)} e^{m \cdot \text{kl}(R_S(h), R(h))} \right] \\
 &\geq \mathbf{E}_{h \sim Q} \ln \left[ \frac{P(h)}{Q(h)} e^{m \cdot \text{kl}(R_S(h), R(h))} \right] \\
 &= m \mathbf{E}_{h \sim Q} \text{kl}(R_S(h), R(h)) - \text{KL}(Q \| P) \\
 &\geq m \cdot \text{kl}(\mathbf{E}_{h \sim Q} R_S(h), \mathbf{E}_{h \sim Q} R(h)) - \text{KL}(Q \| P) \\
 &= m \cdot \text{kl}(R_S(G_Q), R(G_Q)) - \text{KL}(Q \| P) .
 \end{aligned}$$

### Aligned posterior theorem

Here, we do the same operation for “free” (proof on next slide):

$$\begin{aligned}
 & \ln \left[ \mathbf{E}_{h \sim P} e^{m \cdot \text{kl}(R_S(h) \| R(h))} \right] \\
 &= \ln \left[ \mathbf{E}_{h \sim Q} e^{m \cdot \text{kl}(R_S(h) \| R(h))} \right] \\
 &\geq \mathbf{E}_{h \sim Q} \ln \left[ e^{m \cdot \text{kl}(R_S(h), R(h))} \right] \\
 &= m \mathbf{E}_{h \sim Q} \text{kl}(R_S(h), R(h)) \\
 &\geq m \cdot \text{kl}(\mathbf{E}_{h \sim Q} R_S(h), \mathbf{E}_{h \sim Q} R(h)) \\
 &= m \cdot \text{kl}(R_S(G_Q), R(G_Q)) .
 \end{aligned}$$

The two “ $\geq$ ” come from Jensen’s inequality:  $\mathbf{E} f(X) \geq f(\mathbf{E} X)$  for convex  $f$ .

First, note that as we have  $h \in \mathcal{H}^S \Rightarrow -h \in \mathcal{H}^S$  :

$$\mathbf{E}_{h \sim P} e^{m \cdot \text{kl}(R_S(h) \| R(h))} = \int_{h \in \mathcal{H}} P(h) e^{m \cdot \text{kl}(R_S(h) \| R(h))} = \int_{h \in \mathcal{H}} P(-h) e^{m \cdot \text{kl}(R_S(-h) \| R(-h))}.$$

Then,

$$\begin{aligned} 2 \mathbf{E}_{h \sim P} e^{m \cdot \text{kl}(R_S(h) \| R(h))} &= \int_{h \in \mathcal{H}} P(h) e^{m \cdot \text{kl}(R_S(h) \| R(h))} + \int_{h \in \mathcal{H}} P(-h) e^{m \cdot \text{kl}(R_S(-h) \| R(-h))} \\ &= \int_{h \in \mathcal{H}} P(h) e^{m \cdot \text{kl}(R_S(h) \| R(h))} + \int_{h \in \mathcal{H}} P(-h) e^{m \cdot \text{kl}(1 - R_S(h) \| 1 - R(h))} \\ &= \int_{h \in \mathcal{H}} (P(h) + P(-h)) e^{m \cdot \text{kl}(R_S(h) \| R(h))} \\ &= \int_{h \in \mathcal{H}} (Q(h) + Q(-h)) e^{m \cdot \text{kl}(R_S(h) \| R(h))} \\ &= \int_{h \in \mathcal{H}} Q(h) e^{m \cdot \text{kl}(R_S(h) \| R(h))} + \int_{h \in \mathcal{H}} Q(-h) e^{m \cdot \text{kl}(R_S(-h) \| R(-h))} \\ &= 2 \mathbf{E}_{h \sim Q} e^{m \cdot \text{kl}(R_S(h) \| R(h))}. \end{aligned}$$

In this lecture, we will :

- Review quickly the Sample-Compress theory
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers (the Sc-SVM)**
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- Present some experimental **results**
- and Conclude...

# Let's design two learning algorithms

The task of the algorithms is to find a vector  $\mathbf{w} = (w_0, w_1, \dots, w_m)$ ,

$$\begin{aligned}w_0 &\stackrel{\text{def}}{=} w_{\langle \rangle} = Q(h_{\langle \rangle}^{(\sigma,+)}) - Q(h_{\langle \rangle}^{(\sigma,-)}) \\w_i &\stackrel{\text{def}}{=} w_{\langle i \rangle} = Q(h_{\langle i \rangle}^{(\sigma,+)}) - Q(h_{\langle i \rangle}^{(\sigma,-)}) \\|w_j| &\leq \frac{1}{m+1} \quad \forall j \in \{0, \dots, m\}\end{aligned}$$

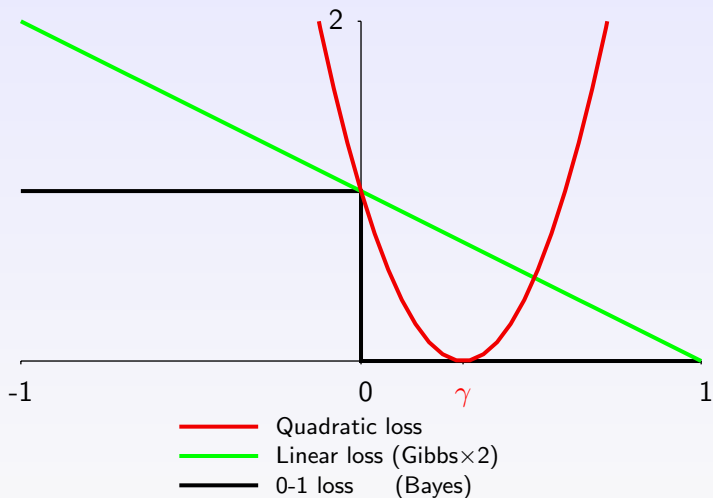
The empirical margin  $\widehat{M}_Q$  will now be defined by

$$\widehat{M}_Q(\mathbf{x}_j, y_j) = \sum_{k=0}^m y_j w_k \widehat{G}(\mathbf{x}_k, \mathbf{x}_j) = y_j \mathbf{w} \widehat{\mathbf{G}}(\mathbf{x}_j)$$

where

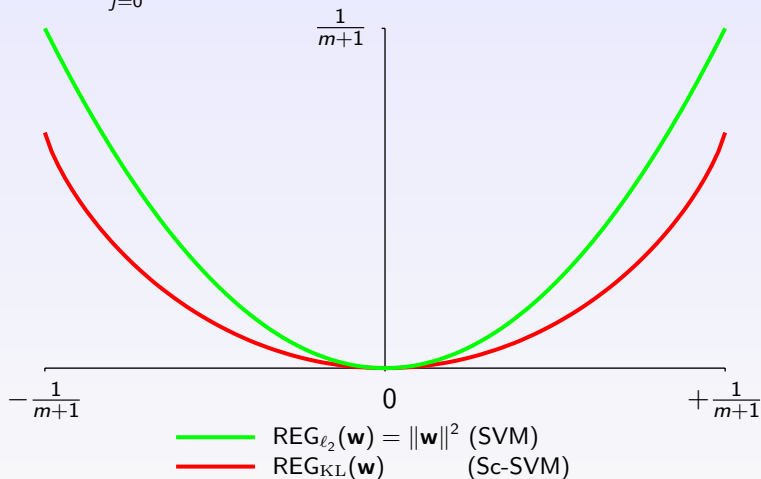
$$\widehat{G}(\mathbf{x}_j, \mathbf{x}_l) \stackrel{\text{def}}{=} \begin{cases} k(\mathbf{x}_j, \mathbf{x}_l) & \forall j \in \{1, \dots, m\} \text{ and } j \neq l \\ 1 & \forall j \in \{1, \dots, m\} \text{ and } j = l \\ 1 & \text{for } j = 0 \end{cases}$$
$$\widehat{\mathbf{G}}(\mathbf{x}_l) \stackrel{\text{def}}{=} (\widehat{G}(\mathbf{x}_0, \mathbf{x}_l), \dots, \widehat{G}(\mathbf{x}_m, \mathbf{x}_l)).$$

Remember that we minimize the **quadratic loss**  $\zeta_\gamma(\alpha) = \left(1 - \frac{1}{\gamma}\alpha\right)^2$ , where  $\alpha$  is the **margin** and  $\gamma$  is the **minimum** of the parabola.



The KL between an uniform prior and the posterior associated to  $\mathbf{w}$  is

$$\text{REG}_{\text{KL}}(\mathbf{w}) = \frac{1}{2} \sum_{j=0}^m \left[ (c + w_j) \ln(c + w_j) + (c - w_j) \ln(c - w_j) \right] - \ln c \quad \text{with } c = \frac{1}{m+1}$$



## Algorithm with $KL$

Find  $\mathbf{w}$  that minimizes  $f(\mathbf{w}) \stackrel{\text{def}}{=} C \cdot \sum_{j=0}^m \zeta_{\gamma} \left( y_j \mathbf{w} \hat{\mathbf{G}}(\mathbf{x}_j) \right) + \text{REG}_{\text{KL}}(\mathbf{w})$

Parameters to tune :

- $C$ , the trade-off between the two terms to minimize
- $\gamma$ , the minimum of the quadratic risk
- Kernel parameter(s), if any

## Algorithm without $KL$

Find  $\mathbf{w}$  that minimizes  $f(\mathbf{w}) \stackrel{\text{def}}{=} \sum_{j=0}^m \zeta_{\gamma} \left( y_j \mathbf{w} \hat{\mathbf{G}}(\mathbf{x}_j) \right)$

Parameters to tune :

- $\gamma$ , the minimum of the quadratic risk
- Kernel parameter(s), if any



## Optimization procedure

Both objective functions are **convex**. Starting from  $\mathbf{w} = \mathbf{0}$ , we optimize  $f(\mathbf{w})$  **coordinate by coordinate**:

- Choose at random  $i \in \{0, \dots, m\}$
- Update  $w_i \leftarrow w_i + \delta$  in order to minimize  $f(\mathbf{w})$
- If  $w_i > \frac{1}{m+1}$  then  $w_i \leftarrow \frac{1}{m+1}$
- If  $w_i < \frac{-1}{m+1}$  then  $w_i \leftarrow \frac{-1}{m+1}$
- Repeat until convergence

Let  $\mathbf{w}_\delta$  be the weight vector obtained after an update  $w_i \leftarrow w_i + \delta$ .

Then, the optimal  $\delta$  is obtain when

$$\frac{\partial f(\mathbf{w}_\delta)}{\partial \delta} = 0$$

## Algorithm with $KL$

Find  $\mathbf{w}$  that minimizes  $f(\mathbf{w}) \stackrel{\text{def}}{=} C \cdot \sum_{j=0}^m \zeta_{\gamma} \left( y_j \mathbf{w} \widehat{\mathbf{G}}(\mathbf{x}_j) \right) + \text{REG}_{KL}(\mathbf{w})$

$$\frac{df(\mathbf{w}_{\delta})}{d\delta} = \frac{2C}{\gamma^2 m} \left( \delta \sum_{j=1}^m \widehat{G}^2(\mathbf{x}_i, \mathbf{x}_j) + \sum_{j=1}^m \widehat{G}(\mathbf{x}_i, \mathbf{x}_j) D_{\mathbf{w}}(j) \right) + \frac{1}{2} \ln \left[ \frac{\frac{1}{m+1} + w_i + \delta}{\frac{1}{m+1} - w_i - \delta} \right]$$

where  $D_{\mathbf{w}}(j) \stackrel{\text{def}}{=} \mathbf{w} \cdot \widehat{\mathbf{G}}(\mathbf{x}_j) - \gamma y_j$ .

We find  $\delta^*$  such as  $\frac{df(\mathbf{w}_{\delta^*})}{d\delta} = 0$  using an iterative **root-finding method**.

## Algorithm without $KL$

Find  $\mathbf{w}$  that minimizes  $f(\mathbf{w}) \stackrel{\text{def}}{=} \sum_{j=0}^m \zeta_{\gamma} \left( y_j \mathbf{w} \widehat{\mathbf{G}}(\mathbf{x}_j) \right)$

$$\frac{df(\mathbf{w}_{\delta})}{d\delta} = \frac{2}{\gamma^2 m} \left( \delta \sum_{j=1}^m \widehat{G}^2(\mathbf{x}_i, \mathbf{x}_j) + \sum_{j=1}^m \widehat{G}(\mathbf{x}_i, \mathbf{x}_j) D_{\mathbf{w}}(j) \right)$$

We find  $\delta^*$  such as  $\frac{df(\mathbf{w}_{\delta^*})}{d\delta} = 0$  computing **directly**  $\delta^* = \frac{-\delta \sum_{j=1}^m \widehat{G}^2(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{j=1}^m \widehat{G}(\mathbf{x}_i, \mathbf{x}_j) D_{\mathbf{w}}(j)}$

In this lecture, we will :

- Review quickly the Sample-Compress theory
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers** (the Sc-SVM)
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- **Present some experimental results**
- and Conclude...

# Experimental results (RBF kernel, 10-folds CV)

Dataset	T	S	n	Classic SVM	SC-SVM (with KL)	SC-SVM (w/o KL)
Usvotes	200	235	16	0.065	<b>0.060</b>	<b>0.060</b>
Liver	175	170	6	<b>0.303</b>	0.371	<b>0.303</b>
Credit-A	300	353	15	0.187	0.170	<b>0.150</b>
Glass	107	107	9	0.159	<b>0.131</b>	0.178
Haberman	150	144	3	<b>0.273</b>	0.287	0.287
Heart	147	150	13	0.184	<b>0.163</b>	0.190
sonar	104	104	60	0.183	0.144	<b>0.135</b>
BreastCancer	340	343	9	0.038	<b>0.035</b>	<b>0.035</b>
Tic-tac-toe	479	479	9	0.023	<b>0.015</b>	<b>0.015</b>
Ionosphere	175	176	34	0.051	<b>0.029</b>	<b>0.029</b>
Wdbc	284	285	30	0.070	0.092	<b>0.067</b>
MNIST:0vs8	1916	500	784	0.005	<b>0.004</b>	<b>0.004</b>
MNIST:1vs7	1922	500	784	0.012	<b>0.008</b>	0.010
MNIST:1vs8	1936	500	784	0.013	<b>0.011</b>	<b>0.011</b>
MNIST:2vs3	1905	500	784	0.023	<b>0.016</b>	0.018
Letter:AB	1055	500	16	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>
Letter:DO	1058	500	16	0.013	<b>0.009</b>	<b>0.009</b>
Letter:OQ	1036	500	16	<b>0.014</b>	0.017	0.017
Adult	10000	1809	14	0.160	<b>0.157</b>	<b>0.157</b>
Mushroom	4062	4062	22	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
Waveform	4000	4000	21	<b>0.068</b>	0.069	<b>0.068</b>
Ringnorm	3700	3700	20	0.015	0.016	<b>0.012</b>

In this lecture, we will :

- Review quickly the Sample-Compress theory
- See how we can describe a SVM as a **Majority Vote of Sample-Compressed classifiers** (the Sc-SVM)
- Use the **PAC-Bayes** theory to **upper-bound** the risk of our Sc-SVM
- Design a **learning algorithm** to minimise this PAC-Bayes bound
- Present some experimental **results**
- **and Conclude...**

# Future works

Two future research ideas (among others) :

- Experimentations with **undefined similiraty measures** (non-PSD kernels)
- Consider a majority vote of sc-classifiers of **maximum size**  $> 1$   
⇒ More general than the SVM



Image: <http://www.mositronic.com/>

Any Questions ?