

# Structural Risk Minimization

Empirical misclassification in Binary classification:

$$\begin{aligned}L_m(\text{sign}(h)) &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i \cdot \text{sign}(h)(x_i) < 0} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i h(x_i) < 0}\end{aligned}$$

Risk convexification: Replace  $z \mapsto \mathbb{1}_{z < 0}$   
by some convex function  $z \mapsto \ell(z)$

Remark (Logistic Regression) In some cases,

risk convexification can be interpreted as maximum

likelihood estimation of some Generalized Lin. Model

Consider a binary response  $Y \in \{0,1\}$ . In logistic regression, the log-likelihood ratio is given by

$$\log \frac{\mathbb{P}_2(Y=1 | X=x)}{\mathbb{P}_2(Y=0 | X=x)} = \beta_0 + \beta^T x$$

that yields

$$\underbrace{\mathbb{P}_2(Y=1 | X=x)}_{\mathbb{E}[Y | X=x]} = \frac{e^{\beta_0 + \beta^T x}}{1 + e^{\beta_0 + \beta^T x}}$$

The negative log-likelihood is:

$$-\frac{1}{n} \sum_{i=1}^n \left\{ y_i \log \mathbb{P}_2[Y=1 | x_i] + (1-y_i) \log \mathbb{P}_2[Y=0 | x_i] \right\}$$

In ML, it is more common to code

$$Y \in \{-1, +1\}$$

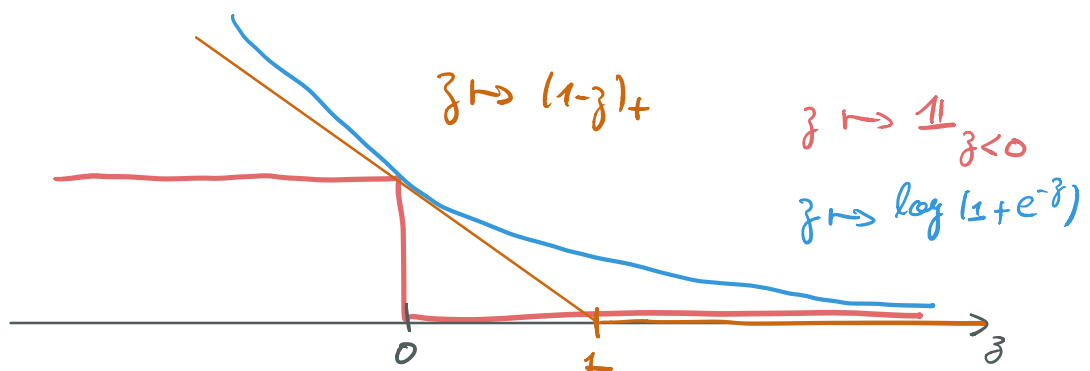
and, in this case, the negative log-likelihood is

$$\frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i h(x_i)})$$

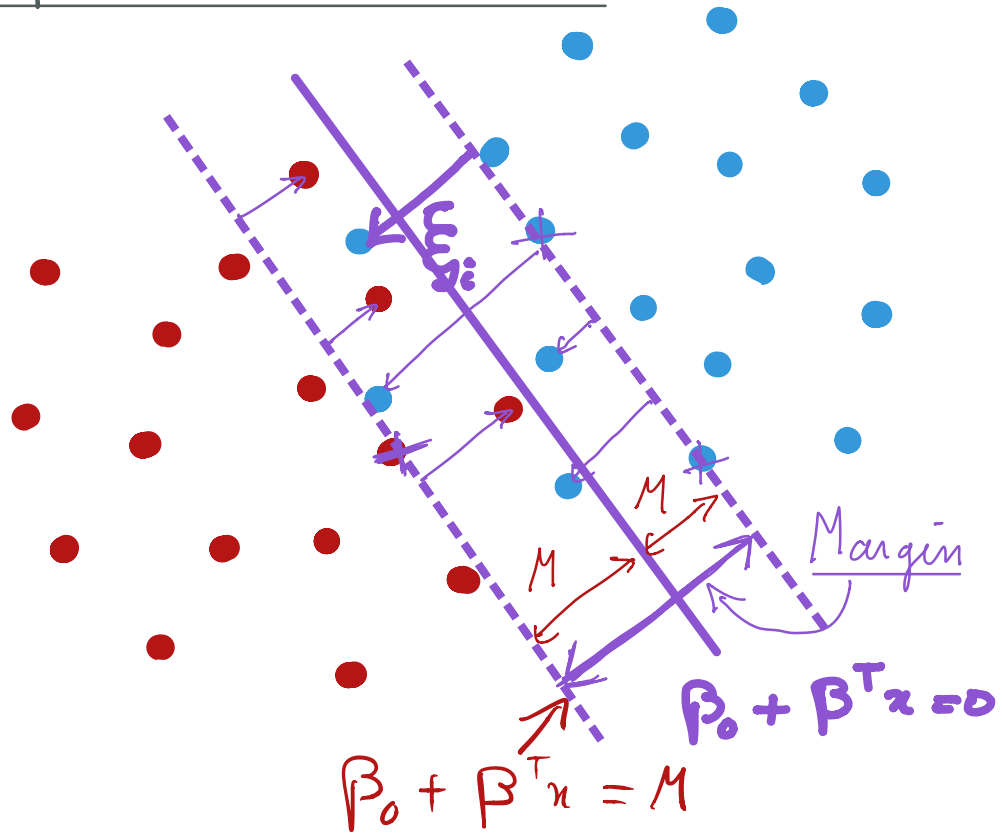
$$\text{with } h(x_i) = \beta_0 + \beta^T x_i$$

### Popular loss functions

- logit loss:  $\log(1 + e^{-z})$
- hinge loss:  $(1 - z)_+ := \max(0, 1 - z)$



# Support Vector Machines



$$l(z) = (1-z)_+$$

$$M = \frac{1}{\|\beta\|_2}$$

Maximize  
 $\beta_0, \beta, \xi_i$

$M$

subject to  $y_i (\beta_0 + \beta^T x_i) \geq M (1 - \xi_i)$

$$\text{and} \left| \begin{array}{l} \xi_i \geq 0 \\ \sum_{i=1}^n \xi_i \leq C \\ \|\beta\|_2 = 1 \end{array} \right.$$

This program is equivalent to:

$$\min_{\beta_0, \beta} \left\{ \frac{1}{n} \sum_{i=1}^n \underbrace{(1 - y_i h(x_i))}_{{= \xi_i}} + \lambda \|\beta\|_2^2 \right\}$$

$$\text{where } h(x) = \beta_0 + \beta^T x$$

solution is given by :

$$\hat{\beta} = \sum_{i=1}^m \hat{\alpha}_i y_i x_i$$

$$\hat{\beta}_0 = \sum_{i=1}^m \hat{\alpha}_i y_i$$

$$\hat{h}(x) = \sum_{i=1}^m \hat{\alpha}_i y_i \left\langle \begin{pmatrix} x_i \\ 1 \end{pmatrix}, \begin{pmatrix} x \\ 1 \end{pmatrix} \right\rangle$$

$$= \hat{\beta}_0 + \hat{\beta}^T x$$

$$x = \begin{pmatrix} x \\ 1 \end{pmatrix} \quad x_i = \begin{pmatrix} x_i \\ 1 \end{pmatrix} :$$

$$\hat{h}(x) = \sum_{i=1}^m \hat{\alpha}_i y_i \langle x_i, x \rangle$$

$$\uparrow$$
$$k(x_i, x)$$

## Reproducing Kernel Hilbert Spaces

- In some ML problems, solutions are given by product of matrices:

$$X_m^T X_m, \quad X_m Y_m \dots$$

where  $X_m = [x_1 \dots x_m]$

$$Y_m = [y_1 \dots y_m]$$

Idea:  $(X_m^T X_m)_{ij} = \langle x_i, x_j \rangle$

$$(X_m Y_m)_i = y_i x_i$$

$$\langle x_i, x_j \rangle \longleftarrow k(x_i, x_j)$$

$$\langle y_i, x_i, x \rangle \longleftarrow y_i k(x_i, x)$$

RKHS:  $\forall n \geq 2, \forall x_1, x_2, \dots, x_n \in \mathcal{X}$

$$K = \left[ k(x_i, x_j) \right]_{1 \leq i, j \leq n} \succeq 0$$

$k(x, y)$  symmetric:

$$k(x, y) = k(y, x)$$

ex:  $k(x, y) = \langle x, y \rangle$



# SVM $\rightarrow$ RKHS framework

---

$$\cdot \frac{1}{m} \sum_{i=1}^m (1 - y_i h(x_i))_+ + \lambda \|\hat{\beta}\|_2^2$$

$$\begin{aligned} \cdot \hat{\beta} &= \sum_{i=1}^m \hat{\alpha}_i y_i \begin{pmatrix} x_i \\ 1 \end{pmatrix} \\ &= \sum_{i=1}^m \hat{c}_i x_i \quad h(x) = x^T [Kc] \end{aligned}$$

$$\cdot \|\hat{\beta}\|_2^2 = \sum_{i,j} \hat{c}_i \hat{c}_j \langle x_i, x_j \rangle$$

$$= c^T K c \quad \text{with } K = (\langle x_i, x_j \rangle)_{i,j}$$

$$\cdot h(x_i) = \langle \hat{\beta}, x_i \rangle = \sum_j \hat{c}_j \langle x_j, x_i \rangle$$

$$= [Kc]_i$$

$$\text{With } \left\{ \begin{array}{l} K = [\langle x_i, x_j \rangle]_{i,j} \\ \beta = \sum c_i \begin{pmatrix} \gamma_i \\ 1 \end{pmatrix} \end{array} \right.$$

$$\left\{ \frac{1}{n} \sum_{i=1}^n (1 - y_i [Kc]_i)_+ + c^T K c \right\}$$

Kernel SVM

Instead of  $K = [\langle x_i, x_j \rangle]_{i,j} = \Phi_n^T \Phi_n$

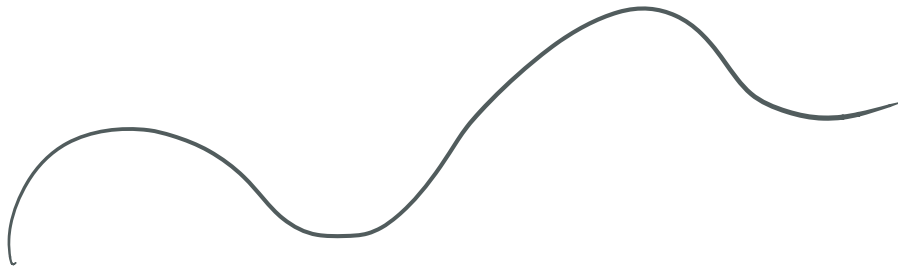
We can use any RKHS matrix  $K$ .

$$K = [k(x_i, x_j)]_{ij}$$

with  $k$  "good" kernel

$$h(x) = \sum_{i=1}^m c_i k(x_i, x)$$

$$\{x : h(x) = 0\} = \left\{x : \sum_{i=1}^m c_i k(x_i, x) = 0\right\}$$



We look at the kernelized SVM:

$$\hat{h}(\cdot) = \sum_{j=1}^m \hat{\beta}_j k(x_j, \cdot)$$

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^m (1 - y_i h(x_i))_+ + \lambda \sum_{i,j} \beta_i \beta_j k(x_i, x_j) \right\}$$

Prop:  $\hat{\beta}$  is such that:

$$\hat{\beta}_i = 0$$

if

$$y_i \hat{h}(x_i) > 1$$

$$\hat{\beta}_i = \frac{y_i}{2\lambda m}$$

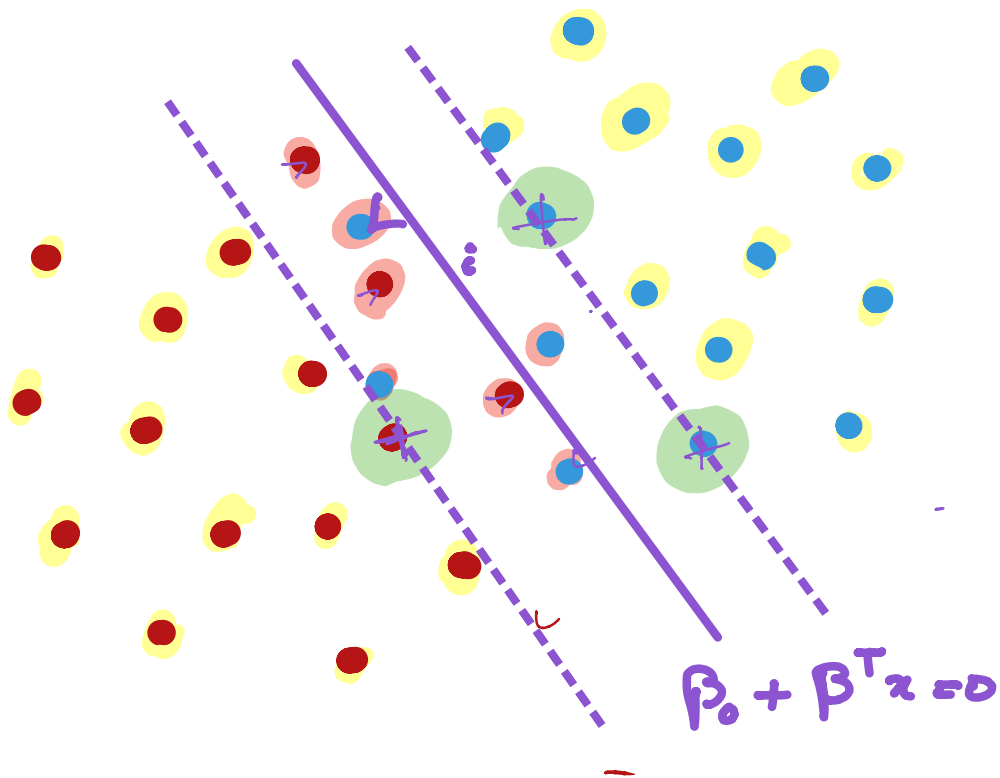
if

$$y_i \hat{h}(x_i) < 1$$

$\hat{\beta}_i$  is such that

$$0 \leq y_i \hat{\beta}_i \leq \frac{1}{2\lambda n}$$

if  $y_i \hat{h}(x_i) = 1$



Proof: Slack variables:

$$\xi_i = \left( 1 - \gamma_i \underbrace{[K\beta]_i}_{h(x_i)} \right)_+$$

Program  $\Leftrightarrow$

$$(\hat{\beta}, \hat{\xi}) \in \arg \min_{\beta, \xi} \left\{ \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \beta^T K \beta \right\}$$

$$\begin{cases} \xi_i \geq 1 - \gamma_i [K\beta]_i \\ \xi_i \geq 0 \end{cases}$$

KKT conditions for the Lagrangian:

$$\mathcal{L}(\beta, \xi, \alpha, \gamma) = \left\{ \frac{1}{n} \sum_i \xi_i + \lambda \beta^T K \beta \right.$$

$$- \sum_{i=1}^n (\alpha_i (\sum_{j=1}^n \xi_j - 1 + Y_i [K\hat{\beta}]_i) + \delta_i \xi_i) \Bigg\}$$

First order:  $2\lambda [K\hat{\beta}]_j = \sum_i K_{ij} \alpha_i Y_i$

and  $\alpha_j + \delta_j = \frac{1}{n}$

Slackness:  $\min(\alpha_i, \sum_{j=1}^n \xi_j - 1 + Y_i [K\hat{\beta}]_i) = 0$   
 $\min(\delta_i, \sum_{j=1}^n \xi_j) = 0$