

PH 240C Supervised Learning (2): Empirical Risk Minimization and Kernel Methods

Jingshen Wang

September 19, 2021

1 Lagrange duality

To solve the problem in Figure 1, we need to work with kernel SVM which is motivated by the dual problem of the optimization problem. Moreover, the dual problem will also allow us to derive an efficient algorithm to solve SVM better than generic QP software.

Review of Lagrange with equality constraints Instead of considering an inequality constraints in the optimization problem for SVM, consider a problem of a simpler form with equality constraint:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} \quad & f(\mathbf{w}) \\ \text{s.t.} \quad & h_i(\mathbf{w}) = 0, \quad i = 1, \dots, n. \end{aligned} \tag{1}$$

We define the **Lagrangian** $L : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$ associated with the problem (1) as

$$L(\mathbf{w}, \boldsymbol{\lambda}) = f(\mathbf{w}) + \sum_{i=1}^n \lambda_i h_i(\mathbf{w}).$$

We refer to λ_i as the Lagrange multiplier associated with the i th equality constraint $h_i(\mathbf{w}) = 0$. The vector $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)'$ is called the **dual variable** or the **Lagrange multiplier vector** associated with the problem (1). We would then find and set $L(\mathbf{w}, \boldsymbol{\lambda})$'s partial derivatives to zero:

$$\frac{\partial L(\mathbf{w}, \boldsymbol{\lambda})}{\partial w_j} = 0, \quad \frac{\partial L(\mathbf{w}, \boldsymbol{\lambda})}{\partial \lambda_i} = 0, \quad j = 1, \dots, d, \quad i = 1, \dots, n,$$

and solve for \mathbf{w} and $\boldsymbol{\lambda}$.

Lagrange with inequality and equality constraints—primal and dual problems Now we generalize this to constrained optimization problems in which we may have inequality as well as equality constraints:

$$\begin{aligned} \text{Primal problem:} \quad & \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ & \text{s.t.} \quad g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, k, \\ & \quad \quad h_i(\mathbf{w}) = 0, \quad i = 1, \dots, n. \end{aligned} \tag{2}$$

We refer to the above problem as the **primal problem**, and the Lagrangian associated with the primal problem is

$$L(w, \lambda, \nu) = f(w) + \sum_{i=1}^k \nu_i g_i(w) + \sum_{i=1}^n \lambda_i h_i(w).$$

Here, $\nu_i \geq 0$ and λ_i are the Lagrange multipliers. We can show that the solution of the primal problem (2) is also the solution to the following problem (see footnote for proof)¹:

$$\min_W L_{\text{primal}}(w) = \min_W \max_{\nu \geq 0, \lambda} L(w, \lambda, \nu).$$

We also define the optimal value of the objective function to be $p^* = \min_W L_{\text{primal}}(w)$.

Next, we define the **Lagrange dual function** (or just dual function) as

$$L_{\text{dual}}(\nu, \lambda) = \min_W L(w, \lambda, \nu).$$

When the Lagrange is unbounded below in w , the dual function takes on the value $-\infty$. The **dual problem** is then defined as

$$\text{Dual problem: } \max_{\nu \geq 0, \lambda} L_{\text{dual}}(w) = \max_{\nu \geq 0, \lambda} \min_W L(w, \lambda, \nu).$$

This is exactly the same as our primal problem shown above, except that the order of the “max” and the “min” are now exchanged. We again define the optimal value of the dual problem’s objective to be $d^* = \max_W L_{\text{dual}}(w)$.

How are the primal and dual problems related? It is easy to see that

$$d^* = \max_W \min_{\nu \geq 0, \lambda} L(w, \lambda, \nu) \leq \min_{\nu \geq 0, \lambda} \max_W L(w, \lambda, \nu) = p^*.$$

Under certain conditions ², we can show that the dual problem and the primal problem have no gap (no

1

Proof. To see this, consider the quantity

$$L_{\text{primal}}(w) = \max_{\nu \geq 0, \lambda} L(w, \lambda, \nu).$$

Then, for a given w : (a) If w violates any of the primal constraints, then we can verify that

$$L_{\text{primal}}(w) = \max_{\nu \geq 0, \lambda} \left(f(w) + \sum_{i=1}^k \nu_i \underbrace{g_i(w)}_{\leq 0} + \sum_{i=1}^n \lambda_i h_i(w) \right) = +\infty.$$

Conversely, if the constraints are indeed satisfied for a particular value of w , then $L_{\text{primal}}(w) = f(w)$. Hence, we conclude

$$L_{\text{primal}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ +\infty & \text{otherwise.} \end{cases}$$

Thus, $L_{\text{primal}}(w)$ takes the same value as the objective function in the problem (2) for all w that satisfies the primal constraints, and is positive infinity if the constraints are violated. We conclude

$$\min_W L_{\text{primal}}(w) = \min_W \max_{\nu \geq 0, \lambda} L(w, \lambda, \nu).$$

Therefore, the solution to the primal problem is also the solution to the unconstrained problem $\min_W \max_{\nu \geq 0, \lambda} L(w, \lambda, \nu)$. □

²One set of sufficient conditions: Suppose f and g_i ’s are convex function, and h_i ’s are affine. Further suppose that the

duality gap): $d^* = p^*$. This suggests there must exist $(\mathbf{w}^*, \nu^*, \lambda^*)$ that satisfies the *Karush-Kuhn-Tucker (KKT)* conditions, which are defined as

$$\begin{aligned}\frac{\partial L(\mathbf{w}^*, \nu^*, \lambda^*)}{\partial \mathbf{w}_i} &= 0, \quad i = 1, \dots, p, \\ \frac{\partial L(\mathbf{w}^*, \nu^*, \lambda^*)}{\partial \lambda_i} &= 0, \quad i = 1, \dots, n \\ \nu_i g_i(\mathbf{w}) &= 0, \quad i = 1, \dots, k \\ g_i(\mathbf{w}) &\leq 0, \quad i = 1, \dots, k \\ \nu_i &\geq 0, \quad i = 1, \dots, k.\end{aligned}$$

2 SVM-continued

2.1 Dual problem of SVM

Let's now circle back to the (primal) SVM optimization problem defined in last class:

$$\begin{aligned}\min_{a, \mathbf{w}} \quad & \frac{1}{2} \mathbf{w}' \mathbf{w} \\ \text{s.t.} \quad & -Y_i(a + \mathbf{w}' X_i) + 1 \leq 0, \quad i = 1, \dots, n.\end{aligned}\tag{3}$$

Borrowing the notation from the previous section, the inequality constraints are

$$g_i(\mathbf{w}) = -Y_i(a + \mathbf{w}' X_i) + 1 \leq 0.$$

Since there is no equality constraint in the problem (3), the Lagrangian for our optimization problem is

$$L(\mathbf{w}, a, \nu) = \frac{1}{2} \mathbf{w}' \mathbf{w} - \sum_{i=1}^n \nu_i [Y_i(a + \mathbf{w}' X_i) - 1].\tag{4}$$

By checking the conditions listed in the previous section, there is no duality gap between the primal problem (3) and its dual. To get the dual problem of (3), we follow the steps introduced in the last section:

1. Find the Lagrange dual function associated with (3):

$$L_{\text{dual}}(\nu) = \min_{a, \mathbf{w}} L(a, \mathbf{w}, \nu).$$

To do so, we take derivative with respect to \mathbf{w} and a :

$$\frac{\partial L(a, \mathbf{w}, \nu)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \nu_i Y_i X_i = 0, \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^n \nu_i Y_i X_i,\tag{5}$$

$$\frac{\partial L(a, \mathbf{w}, \nu)}{\partial a} = \sum_{i=1}^n \nu_i Y_i = 0, \quad \Rightarrow \quad \sum_{i=1}^n \nu_i Y_i = 0.\tag{6}$$

constraint g_i are (strictly) feasible—meaning that there exists some \mathbf{w} such that $g_i(\mathbf{w}) < 0$ for all i . Note there are many different sets of sufficient conditions can be found in the literature.

2. We now take the definition of w in (5) and plug it back into the equation (4):

$$\begin{aligned} L_{\text{dual}}(\nu) &= \min_{a, w} L(a, w, \nu) = \sum_{i=1}^n \nu_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n Y_i Y_j \nu_i \nu_j X_i' X_j - a \sum_{i=1}^n \nu_i Y_i \\ &= \sum_{i=1}^n \nu_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n Y_i Y_j \nu_i \nu_j X_i' X_j. \end{aligned}$$

3. **SVM dual problem:** Dual problem associated with the primal problem (3) is thus

$$\begin{aligned} \max \quad & L_{\text{dual}}(\nu) = \sum_{i=1}^n \nu_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n Y_i Y_j \nu_i \nu_j X_i' X_j \\ \text{s.t.} \quad & \nu_i \geq 0, \quad i = 1, \dots, n, \\ & \sum_{i=1}^n \nu_i Y_i = 0. \end{aligned} \tag{7}$$

4. **Primal and dual relationship:** Suppose we can solve the dual problem (for now) and the optimal solution is ν^* , then we can obtain the primal solution from (5)

$$w^* = \sum_{i=1}^n \nu_i^* Y_i X_i,$$

how about the intercept a^* ?

We want to resolve several issues in SVM by working out the dual formulation:

First, what does $\nu_i^* = 0$ tell us? Recall from the KKT condition, we must have (this specific condition is also called KKT dual complementarity condition)

$$\nu_i [Y_i (a + w' X_i) - 1] = 0.$$

It says that if $\nu_i^* > 0$, then $Y_i (a + w' X_i) - 1 = 0$, meaning that i th data point is a support vector. Whenever $\nu_i^* = 0$, we have the sample pair (Y_i, X_i) does not live on the decision boundary.

Second, when we make prediction at a new data point $x \in \mathbb{R}^d$, we would then calculate $w^{*'} x + a^*$ and then predict

$$y = \begin{cases} 1 & w^{*'} x + a^* \geq 0 \\ -1 & w^{*'} x + a^* < 0. \end{cases}$$

Notice that the quantity $w^{*'} x + a^*$ can be rewritten as

$$\begin{aligned} w^{*'} x + a^* &= \left(\sum_{i=1}^n \nu_i^* Y_i X_i \right)' x + a^* \\ &= \sum_{i=1}^n \nu_i^* Y_i X_i' x + a^* \\ &= a^* + \sum_{i: \nu_i^* \neq 0} \mathbf{1}_{(Y_i=1)} X_i' x - \sum_{i: \nu_i^* \neq 0} \mathbf{1}_{(Y_i=-1)} X_i' x. \end{aligned}$$

Now we can see that decision made from SVM is also quite intuitive: if the new data point x is more “similar” to the support vectors labelled as one—hence there is a higher “chance” that $\sum_{i:\nu_i \neq 0} \mathbf{1}_{(Y_i=1)} X'_i x > \sum_{i:\nu_i \neq 0} \mathbf{1}_{(Y_i=-1)} X'_i x$ —thereby it classifies x as 1. Can this motivate you to classify sample in a different way?

By examining the dual form of the optimization problem, we gained significant insight into the structure of the problem, and were also able to write the entire algorithm in terms of only inner products between input feature vectors. In the next section, we will exploit this property to apply the kernels to our classification problem. The resulting algorithm, support vector machines, will be able to efficiently learn in very high dimensional spaces.

2.2 Kernel SVM

From our drug dosage example, we can transform our original input *attributes* to new *features* as in Figure 1.

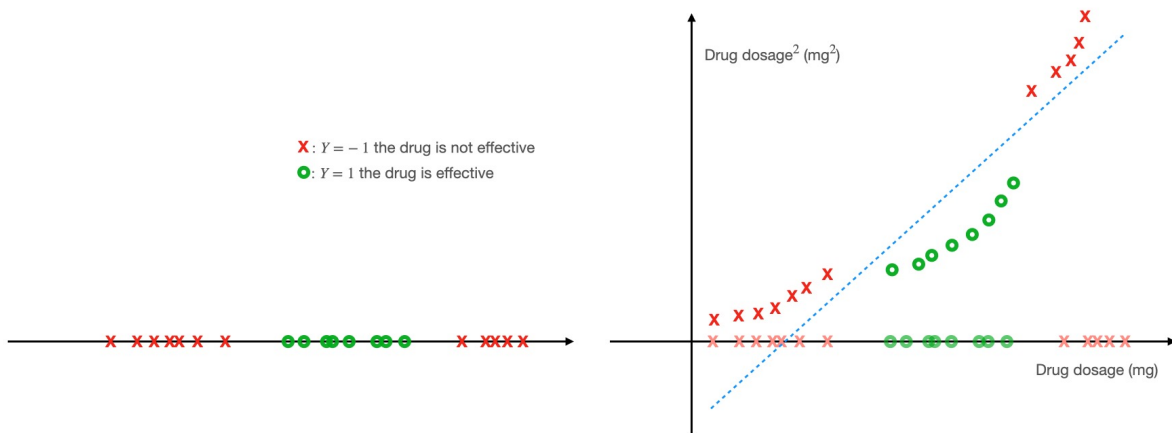


Figure 1: Can SVM be used to predict the drug-effective outcome? Yes!

Let’s formalize this transformation. Let ϕ denote the *feature mapping*, which maps from the attributes to the features. For instance, in our example, we use

$$\phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}.$$

Rather than applying SVMs using the original input attributes x , we may instead want to learn using some features $\phi(x)$. To do so, we simply need to go over our previous algorithm, and replace X_i everywhere in it with $\phi(X_i)$.

Since the SVM dual problem (7) solely works with the inner products $\langle X_i, X_j \rangle$, this means that we would replace all those inner products with $\langle \phi(X_i), \phi(X_j) \rangle$. More specifically, given a feature mapping ϕ , we define its corresponding *Kernel* to be

$$K(X_i, X_j) = \langle \phi(X_i), \phi(X_j) \rangle = \phi(X_i)' \phi(X_j).$$

Then, everywhere we previously had $\langle X_i, X_j \rangle$, we could simply replace it with $K(X_i, X_j)$ and our algorithm would now be learning using the feature mapping ϕ .

Suppose $x, z \in \mathcal{X} \subseteq \mathbb{R}^d$, given a feature mapping ϕ , we could easily compute $K(x, z)$ by finding $\phi(x)$, $\phi(z)$ and taking their inner product. But what's more interesting is that, very often, $K(x, z)$ is very inexpensive to calculate, even though $\phi(x)$ itself may be very expensive to calculate. In such setting, kernel SVM incorporates easy-to-compute kernel function $K(x, z)$, and hence can learn in high-dimensional feature space (given by ϕ) without having to explicitly find or represent $\phi(x)$.

Example 1 (Degree- p polynomials kernel). *For two input attributes defined in the attribute space \mathcal{X} , $x, z \in \mathcal{X} \subseteq \mathbb{R}^d$, we consider a Kernel*

$$K(x, z) = (x'z + c)^p.$$

When $p = 2$, we have

$$\begin{aligned} K(x, z) &= (x'z + c)^2 \\ &= \sum_{k=1}^d \sum_{l=1}^d (x_k x_l)(z_k z_l) + \sum_{k=1}^d (2\sqrt{c}x_k)(2\sqrt{c}z_k) + c^2. \end{aligned}$$

This means if we direct work with the original feature mapping, we have to work with a long feature vector. Take $d = 3$ for example:

$$\phi(x) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_1 x_3 \\ x_2^2 \\ x_2 x_1 \\ x_2 x_3 \\ x_3 x_2 \\ x_3^2 \\ \sqrt{2c}x_1 \\ \sqrt{2c}x_2 \\ \sqrt{2c}x_3 \\ c \end{bmatrix} \in \mathbb{R}^{12}, \quad \text{where } x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^{d=3}.$$

More broadly speaking, the degree- p kernel $K(x, z) = (x'z + c)^p$ corresponds to a feature mapping to an $\binom{p+d}{d}$ -dimensional feature space. Such a feature space is spanned by all monomials of the x_{i_1}, \dots, x_{i_k} that are up to order p - thus this is a space of order d^p . Nevertheless, computing $K(x, z)$ takes only $O(d)$ time. Hence, we never need to explicitly represent feature vector in this very high-dimensional feature space.

Although the above description sounds mathematically complex, the kernel function is essentially measuring the similarity between two input attributes x and z . As long as you come up with some function $K(x, z)$ that you think might be a reasonable measure of how similar x and z are, it would be a reasonable kernel. Another popular choice is the Gaussian Kernel:

Example 2 (Gaussian Kernel). *For two input attributes defined in the attribute space \mathcal{X} , $x, z \in \mathcal{X} \subseteq \mathbb{R}^d$, we consider a Kernel*

$$K(x, z) = \exp\left(-\frac{\|x - z\|_2^2}{2\sigma^2}\right)$$

This Gaussian kernel corresponds to an infinite dimensional feature mapping ϕ . But think about more broadly, given some random function K , how can we tell if it is a valid kernel? In other words, can we tell if there exist a feature mapping ϕ so that $K(x, z) = \phi(x)' \phi(z)$ for all x, z ? This requires us to understand the Mercer's Theorem. If you are interested, you can take a look at Chapter 14.2 in [Murphy \(2012\)](#).

Keep in mind however that the idea of kernels has significantly broader applicability than SVMs. Specifically, if you have any learning algorithm that you can write in terms of only inner products $\langle x, z \rangle$ between input attribute vectors, then by replacing this with $K(x, z)$ where K is a kernel, you can “magically” allow your algorithm to work efficiently in the high dimensional feature space corresponding to K .

Lastly, to end this section on the support vector machines, we summarize its primal and dual problems:

$$\begin{aligned} \text{SVM primal problem: } \quad & \min_{a, \mathbf{w}} \quad \frac{1}{2} \mathbf{w}' \mathbf{w} \\ & \text{s.t.} \quad -Y_i(a + \mathbf{w}' \phi(X_i)) + 1 \leq 0, \quad i = 1, \dots, n. \end{aligned} \tag{8}$$

and the dual problem

$$\begin{aligned} \text{SVM dual problem: } \quad & \max_{\nu} \quad \sum_{i=1}^n \nu_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n Y_i Y_j \nu_i \nu_j K(X_i, X_j) \\ & \text{s.t.} \quad \nu_i \geq 0, \quad i = 1, \dots, n, \\ & \quad \quad \sum_{i=1}^n \nu_i Y_i = 0, \end{aligned} \tag{9}$$

where $K(X_i, X_j) = \langle \phi(X_i), \phi(X_j) \rangle$. The primal and dual solution can be linked by

$$\mathbf{w} = \sum_{i=1}^n \nu_i Y_i \phi(X_i).$$

2.3 A short story of SVM

The very first paper on the support vector machines with kernels is published in 1992; see [Boser et al. \(1992\)](#). You can read a short story on [this webpage](#).

3 Empirical Risk Minimization

Recall in the SVM with soft constraint problem, we work with

$$\begin{aligned} \min_{a, \mathbf{w}} \quad & \mathbf{w}' \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & Y_i(a + \mathbf{w}' X_i) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, \dots, n. \end{aligned}$$

Whenever $C \neq 0$, the optimization problem will always try to minimize ξ_i as much as possible. Then the equation must hold as an equality and we have

$$\xi_i = \begin{cases} 0 & \text{if } Y_i(a + w'X_i) \geq 1 \\ 1 - Y_i(a + w'X_i) & \text{if } Y_i(a + w'X_i) < 1. \end{cases}$$

Therefore, the slack variable ξ_i has a closed form expression:

$$\xi_i = \max\{1 - Y_i(a + w'X_i), 0\}$$

When we plug this closed form into the objective of the SVM primal problem, we obtain the following unconstrained version as loss function:

$$\min_{a, w} C \sum_{i=1}^n \underbrace{\max[1 - Y_i(a + w'X_i), 0]}_{\text{hinge loss}} + \underbrace{w'w}_{l_2 \text{ penalty}}.$$

The hinge loss is intuitive to understand: the classifier h learnt by SVM is parametrized by w , and we make decisions based on

$$h_w(X_i) = \begin{cases} 1 & a + w'X_i \geq 1 \\ -1 & a + w'X_i \leq -1 \end{cases} = 2\mathbf{1}(a + w'X_i \geq 1) - 1 \triangleq 2\mathbf{1}(f_w(X_i) \geq 1) - 1.$$

Whenever the classifier makes a mistake, we have $Y_i \cdot (a + w'X_i) < 0$ —generating a loss. Therefore, we define a hinge loss for the classifier h_w as

$$l(h_w(X_i), Y_i) = \max[1 - Y_i(a + w'X_i), 0] \triangleq \max[1 - Y_i f_w(X_i), 0].$$

Then the SVM primal problem can be understood from a penalized “regression” perspective (suppose we ignore the intercept b for simplicity³)

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n l_{\text{svm}}(h_w(X_i), Y_i) + \lambda \cdot \text{Penalty}(w).$$

In your lab, you should have worked with GLM Lasso (or some general penalty) when we label the outcome as $Y_i \in \{-1, 1\}$ (try this one on your own, the optimization problem we worked out in GLM section is for $Y_i \in \{0, 1\}$):

$$\begin{aligned} \hat{\beta} &= \arg \min_{b \in \mathbb{R}^d} \sum_{i=1}^n \log(1 + \exp(-X_i' b \cdot Y_i)) + \lambda \cdot \text{Penalty}(b), \\ &= \arg \min_{b \in \mathbb{R}^d} \sum_{i=1}^n l_{\text{logit}}(h_b(X_i), Y_i) + \lambda \cdot \text{Penalty}(b), \end{aligned}$$

³This works as once we have solved SVM for some w^* , then optimal value for b^* can be calculated immediately.

where the classifier h_b does the following:

$$h_b(X_i) = \begin{cases} 1 & \text{if } b'X_i \geq 0 \\ -1 & \text{if } b'X_i < 0. \end{cases}$$

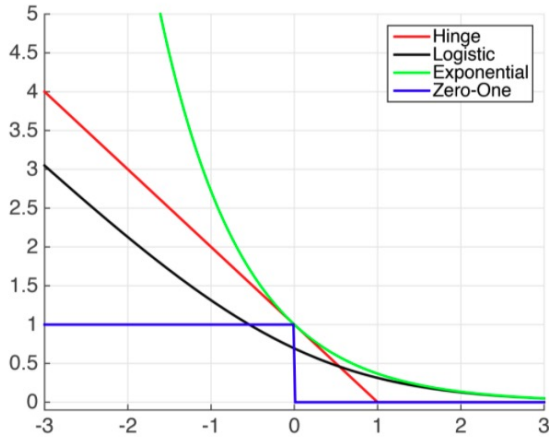


Figure 2: Common Classification Loss Functions on the x-axis: $f(X_i)Y_i$

from a statistical *inference* point of view.

We end this section by defining the general empirical risk minimization problem with a loss function l and some penalty function p :

$$\min_{\mathbf{w}} \underbrace{\frac{1}{n} \sum_{i=1}^n l(h_{\mathbf{w}}(X_i), Y_i)}_{\text{loss}} + \underbrace{\lambda \cdot p(\mathbf{w})}_{\text{penalty}},$$

where the loss function $l(\cdot, \cdot)$ is a function that penalize the difference between Y_i and $h_{\mathbf{w}}(X_i)$, and the penalty function $p(\cdot)$ penalizes the classifier's complexity.

Now, the logistic regression and SVM look oddly similar. Both say that if we misclassify more pairs (X_i, Y_i) , the loss function (either l_{logit} or l_{svm}) is large; the penalty is trying to make the classifier h simple (remember the Occam's razor principle?). In other words, we are trying to find a classifier that (1) makes less mistakes, and (2) as simple as possible. This principle applies in many famous estimators (ridge regression, Lasso, Elastic Net, etc.).

Let's now look into the difference between the SVM and the logistic losses in Figure 2. The exponential loss is of the form (this is a very sensitive/aggressive loss)

$$l_{\text{exponential}}(h(X_i), Y_i) = \exp(-h(X_i)Y_i).$$

Which loss do you like better? Why? Think about it

References

Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, 1992.

Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.