

## Kernel Machines Review / Functional Gradient Descent

Lecturer: Drew Bagnell

Scribe: Abhinav Shrivastava<sup>1</sup>

## 1 The Big Picture

We have seen how to use online convex programming to learn linear functions by optimizing costs of the following form:

$$L(\mathbf{w}) = \sum_i \underbrace{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}_{\text{loss}} + \underbrace{\lambda \|\mathbf{w}\|^2}_{\text{regularization/prior}}$$

We want generalize this to learn over a space of more general functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . The high-level idea is to learn non-linear models using the same gradient-based approach used to learn linear models.

$$L(f) = \sum_i (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|^2$$

Up till now we have only considered functions of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , but we will now extend this to a more general space of functions.

## 2 Review of Kernels

- Ultimately, we wish to learn a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that assigns a meaningful score given a data point. For example, in binary classification, we would like an  $f(\cdot)$  to return positive and negative values, given positive and negative samples, respectively.
- A kernel  $K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  intuitively measures the *correlation* between  $f(\mathbf{x}_i)$  and  $f(\mathbf{x}_j)$ . Considering a matrix  $\mathbf{K}$  with entries  $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ , then matrix  $\mathbf{K}$  must satisfy the properties:

- $\mathbf{K}$  is symmetric ( $K_{ij} = K_{ji}$ )
- $\mathbf{K}$  is positive-definite ( $\forall \mathbf{x} \in \mathbb{R}^n : \mathbf{x} \neq \mathbf{0}, \mathbf{x}^T \mathbf{K} \mathbf{x} > 0$ )

Hence, a valid kernel is the inner product:  $K_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ .

- A function can be considered that is a weighted composition of many kernels centered at various locations  $\mathbf{x}_i$ :

$$f(\cdot) = \sum_{i=1}^Q \alpha_i K(\mathbf{x}_i, \cdot), \quad (1)$$

---

<sup>1</sup>Based on the scribe work of Varun Ramakrishna, Dave Rollinson, Daniel Munoz, Tomas Simon, Jack Singleton and Sergio Valcarcel

where  $Q$  is the number of kernels that compose  $f(\cdot)$  and  $\alpha_i \in \mathbb{R}$  is each kernel's associated weight. All functions  $f(\cdot)$  with kernel  $K$  that satisfy the above properties and can be written in the form of Equation 1 are said to lie in a *Reproducing Kernel Hilbert Space* (RKHS)  $\mathcal{H}_K$ :  $f \in \mathcal{H}_K$

However to do gradient descent on the space of such functions, we need the notion of a distance, norm and an inner product. We formalize this by introducing the Reproducing Kernel Hilbert Space.

### 3 Reproducing Kernel Hilbert Space

The Reproducing Kernel Hilbert Space (RKHS), denoted by  $\mathcal{H}_k$ , is the space of functions  $f(\cdot)$  that can be written as  $\sum_i \alpha_i k(\mathbf{x}_i, \cdot)$ , where  $k(\mathbf{x}_i, \mathbf{x}_j)$  satisfies certain properties described below.

To be able to manipulate objects in this space of functions, we will look at some key properties:

- The **inner product** of  $f, g \in \mathcal{H}_k$  is defined as

$$\langle f, g \rangle \triangleq \sum_i \sum_j \alpha_i \beta_j k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\beta}$$

where  $f(\cdot) = \sum_i \alpha_i k(\mathbf{x}_i, \cdot)$ ,  $g(\cdot) = \sum_j \beta_j k(\mathbf{x}_j, \cdot)$ ,  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  are vectors comprising, respectively,  $\alpha_i$  and  $\beta_i$  components, and  $\mathbf{K}$  is  $n$  by  $m$  (where  $n$  is the number of  $\mathbf{x}_i$  in  $f$ , and  $m$  those in  $g$ ) with  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ .

Note that this will satisfy linearity (in both arguments):

- $\langle \lambda f, g \rangle = \lambda \langle f, g \rangle$
- $\langle f_1 + f_2, g \rangle = \langle f_1, g \rangle + \langle f_2, g \rangle$

With this inner product, the *norm* will be:  $\|f\|^2 = \langle f, f \rangle = \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$ .

- The **reproducing property** is observed by taking the inner-product of a function with a kernel  $\langle f, K(\mathbf{x}^*, \cdot) \rangle$  and functional  $E$ :

$$\begin{aligned} E_{x^*}[f] &= \langle f, K(\mathbf{x}^*, \cdot) \rangle \\ &= \left\langle \sum_{i=1}^Q \alpha_i K(\mathbf{x}_i, \cdot), K(\cdot, \mathbf{x}^*) \right\rangle = \sum_{i=1}^Q \alpha_i \langle K(\mathbf{x}_i, \cdot), K(\cdot, \mathbf{x}^*) \rangle = \sum_{i=1}^Q \alpha_i K(\mathbf{x}_i, \mathbf{x}^*) \\ &= \underbrace{f(\mathbf{x}^*)}_{\text{eval @ } x^*} \end{aligned}$$

An example of a valid kernel for  $\mathbf{x} \in \mathbb{R}^n$  is the inner product:  $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$ . Intuitively, the kernel measures the *correlation* between  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .

A very commonly used kernel is the RBF or Radial Basis Function kernel, which takes the form  $k(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{1}{\gamma} \|\mathbf{x}_i - \mathbf{x}_j\|^2}$ . With this kernel in mind, a function can be considered as a weighted (by

$\alpha_i$ ) composition of bumps (the kernels) centered at the  $Q$  locations  $\mathbf{x}_i$ :

$$f(\cdot) = \sum_{i=1}^Q \alpha_i K(\mathbf{x}_i, \cdot),$$

## 4 Loss Minimization

Again, let us consider our cost function defined over all functions  $f$  in our RKHS, as before our loss is:

$$L(f) = \sum_i (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|^2$$

The purpose of  $\langle f, f \rangle$  is to penalize the complexity of the solution  $f$ . Here it acts like the log of a gaussian prior over functions. Intuitively, the probability can be thought of as being distributed according to  $P(f) = \frac{1}{Z} e^{-\frac{1}{2}\langle f, f \rangle}$  (in practice this expression doesn't work because  $Z$  becomes infinite).

We want to find the best function  $f$  in our RKHS so as to minimize this cost, and we will do this by moving in the direction of the negative gradient:  $f - \alpha \nabla L$ . To do this, we will first have to be able to express the gradient of a function of functions (ie. a *functional* such as  $L[f]$ ).

### 4.1 Functional gradient

A gradient can be thought of as:

- Vector of partial derivatives
- Direction of steepest ascent
- Linear approximation of the function (or functional), ie.  $f(x_0 + \epsilon) = f(x_0) + \epsilon \underbrace{\cdot \nabla f(x_0)}_{\text{gradient}} + O(\epsilon^2)$ .

We will use the third definition. A *functional*  $E : f \rightarrow \mathbb{R}$  is a function of functions  $f \in \mathcal{H}_K$ . As an example let us write the terms of our loss function from above as functionals:

- $E_1[f] = \|f\|^2$
- $E_2[f] = (y - f(\mathbf{x}))^2$
- $E[f] = \sum_i (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|^2$

A functional gradient  $\nabla E[f]$  is defined implicitly as the linear term of the change in a function due to a small perturbation  $\epsilon$  in its input:  $E[f + \epsilon g] = E[f] + \epsilon \langle \nabla E[f], g \rangle + O(\epsilon^2)$ .

Before computing the gradients for these functionals, let us look at a few tools that will help us derive the gradient of the loss functional

## 4.2 Chain rule for functional gradients

Consider *differentiable* functions  $C : \mathbb{R} \rightarrow \mathbb{R}$  that are functions of functionals  $G$ ,  $C(G[f])$ . Our cost function  $L[f]$  from before was such a function, these are precisely the functions that we are interested in minimizing.

The derivative of these functions follows the chain rule:

$$\nabla C(G[f]) = \frac{\partial C(G[f])}{\partial \lambda} \Big|_{G(f)} \nabla G[f] \quad (2)$$

Example: If  $C = (\|f\|^2)^3$ , then  $\nabla C = 3(\|f\|^2)^2(2f)$

## 4.3 Another useful functional gradient

Another useful function that we come across often is the evaluation functional. The *evaluation functional* evaluates  $f$  at the specified  $x$ :  $E_x[f] = f(x)$

- Gradient is  $\nabla E_x = K(x, \cdot)$

$$\begin{aligned} E_x[f + \epsilon g] &= f(x) + \epsilon g(x) + 0 \\ &= f(x) + \epsilon \langle K(x, \cdot), g \rangle + 0 \\ &= E_x[f] + \epsilon \langle \nabla E_x, g \rangle + O(\epsilon^2) \end{aligned}$$

- It is called a *linear functional* due to the lack of a multiplier on perturbation  $\epsilon$ .

## 4.4 Functional gradient of the regularized least squares loss function

- Let's look at the functional gradient of the second term of the loss function:

$$\nabla E[f] = \nabla \|f\|^2 \quad (3)$$

Expanding it out using a Taylor's series type expansion

$$\begin{aligned} E[f + \epsilon g] &= \langle f + \epsilon g, f + \epsilon g \rangle \\ &= \|f\|^2 + 2\langle f, \epsilon g \rangle + \epsilon^2 \|g\|^2 \\ &= \|f\|^2 + \epsilon \langle 2f, g \rangle + O(\epsilon^2) \end{aligned}$$

We observe that

$$\nabla E[f] = \nabla \|f\|^2 = 2f \quad (4)$$

- Now for the first term of the loss function

$$E[f] = \sum_i (y_i - f(\mathbf{x}_i))^2 \quad (5)$$

Using the chain rule we have

$$\nabla E[f] = -2(y_i - f(\mathbf{x}_i)) \nabla(f(x_i)) \quad (6)$$

We observe that  $\nabla(f(x_i))$  is the functional gradient of the evaluation functional. Substituting in the gradient of the evaluation functional as computed in the previous section we have :

$$\nabla E[f] = -2(y_i - f(\mathbf{x}_i)) K(\mathbf{x}_i, \cdot) \quad (7)$$

## 5 Functional gradient descent

- Regularized least squares loss function  $L[f]$

$$\begin{aligned} L[f] &= (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|^2 \\ L[f] &= (y_i - E_{\mathbf{x}_i}[f])^2 + \lambda \|f\|^2 \\ \nabla L[f] &= -2(y_i - f(\mathbf{x}_i))K(\mathbf{x}_i, \cdot) + 2\lambda f \end{aligned}$$

Update rule for the regularized least squares loss function:

$$\begin{aligned} f_{t+1} &\leftarrow f_t - \eta_t \nabla L \\ &\leftarrow f_t - \eta_t (-2(y_t - f_t(\mathbf{x}_t))K(\mathbf{x}_t, \cdot) + 2\lambda f_t) \\ &\leftarrow f_t(1 - 2\eta_t \lambda) + 2\eta_t (y_t - f_t(\mathbf{x}_t))K(\mathbf{x}_t, \cdot) \end{aligned}$$

where  $\eta_t$  is the learning rate at time step  $t$ .

The update rule is equivalent to:

- Adding a kernel  $K(\mathbf{x}_t, \cdot)$  weighted by  $2\eta_t(y_t - f_t(\mathbf{x}_t))$ .
- Shrinking all other weights by  $(1 - 2\eta_t \lambda)$  multiplier.

- SVM loss function  $L(f)$

$$L(f(\mathbf{x}_t), y_t) = \max(0, 1 - y_t f(\mathbf{x}_t)) + \lambda \|f\|^2 \quad (8)$$

The sub-gradient  $\nabla L$  has two cases. One where the prediction is correct by margin = 1, and the other where is not correct by margin = 1 (margin error).

$$\nabla L((x_t), y_t) = \begin{cases} 0 & \text{if } (1 - y_i f(\mathbf{x}_i)) \leq 0 \\ L'(f(\mathbf{x}_t), y_t) f'(\mathbf{x}_t) = -y_t K(\mathbf{x}_t, \cdot) & \text{else margin error} \end{cases} \quad (9)$$

The update rule is equivalent to:

- Adding a kernel  $K(\mathbf{x}_t, \cdot)$  weighted by  $\eta_t y_t$  in case of margin error.
- Shrinking all other weights by  $(1 - 2\eta_t \lambda)$  multiplier.

## 6 Online Kernel Machine

- Initialize the function  $f = 0$ .
- For  $t = 1$  to  $T$ :
  1. Observe some measurement over some set of features  $x_t$
  2. Predict the class using  $f(x_t) = \sum_{i=1}^n \alpha_i K(x_i, x_t)$
  3. Receive loss based on the prediction from  $f(x_t)$  and the true class  $y_t$

$$L(f(x_t), y_t)$$

4. Update  $f$  based on the gradient of the loss function  $L$  and learning rate  $\eta_t$  depending on the chosen algorithm (examples in previous section).

## 6.1 Discussion

- **Representer Theorem** (informally): Given a loss function and regularizer objective with many data points  $\{x_i\}$ , the minimizing solution  $f^*$  can be represented as

$$f^*(\cdot) = \sum_i \alpha_i K(x_i, \cdot) \quad (10)$$

- This algorithm qualitatively corresponds to adding weighted 'bumps' that predicts some value based on the kernel function in each new observation's neighborhood of the feature space in  $x$ . For example: Figure 1 shows an update over 3 points  $\{(x_1, +), (x_2, -), (x_3, +)\}$ . The individual kernels centered at the points are **independently** drawn with colored lines. After 3 updates, the function  $f$  looks like the solid black line.

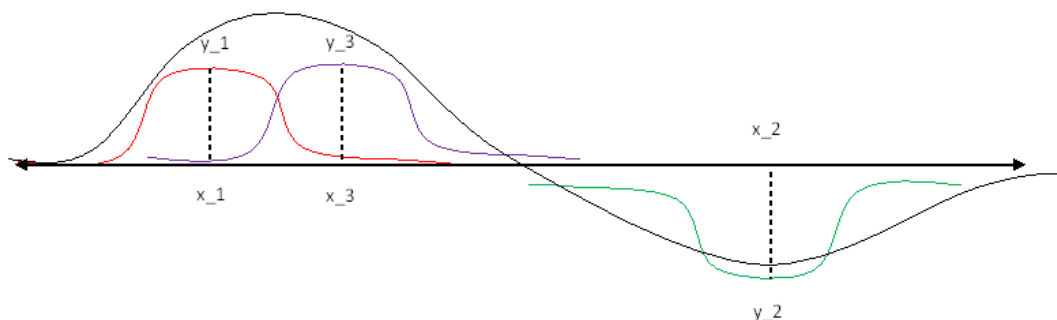


Figure 1: Illustration of function after 3 updates

- Need to perform  $O(T)$  work at each time step. As time progresses and the data set grows, the prediction step will take longer and longer to compute. To shorten this computation time you may want to throw out old data points by weight or age. If interested, there are some papers on that use tricks to find sparse solutions to large-scale problems:

- Rahimi and Recht - *Random Features for Large-Scale Kernel Machines* 2007
- Dekel, Shalev-Shwartz and Singer *The Forgetting: A Kernel-Based Perceptron on a Budget* 2007

- The regret is computed as:

$$Regret = \sum_t (C_t(f_t(x_t)) - C_t(f^*(x_t))) \mid f^* \in H_k$$

The regret bound:

$$Regret = \|\nabla C_t(f)\|_k \cdot \|f^*\|_k \sqrt{T}$$

$\|f^*\|$  is the size of the function.  $\|\nabla C_t(f)\|$  can get as big as  $\alpha^T K \alpha$ .

- The choice and tuning of the kernel and their corresponding bandwidth parameters are what affect the bias-variance tradeoff. These are parameters that need to be tuned in addition to the learning rate  $\eta$  and decay rate  $\lambda$  from the update equations.
  - Often simple kernels work quite well. When approaching a new problem it is usually a good idea start with linear or polynomial kernels. Radial basis functions are another good kernel to try early on. Note that any kernels  $K_1$  and  $K_2$  that satisfy the conditions mentioned in Section 2 can be summed to form a new valid kernel.