

Lecture 1: An Introduction to Supervised Learning

July 2

Lecturer: Lei Wu

Scribe: Shuhai Zhao, Yilei Han

1 Supervised Learning

Some basic terminologies:

- *features*: The set of attributes, often represented as a vector, associated to an example.
- *Hypothesis space*: A set \mathcal{F} of functions mapping features to the set of labels \mathcal{Y} .
- *Loss function*: A function l that measures the difference, or loss, between a predicted label and a true label: $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, for example, $l(y, y') = (y - y')^2$.

Problems in which the training data comprises examples of the input vectors along with their corresponding target vectors are known as *supervised learning* problems. Specifically, in the most general scenario of supervised learning, the distribution D is defined over $\mathcal{X} \times \mathcal{Y}$, and the training data is a labeled sample S drawn i.i.d. according to D :

$$S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)),$$

where $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}, i = 1, 2, \dots, n, y_i = f^*(\mathbf{x}_i) + \epsilon_i$.

Definition 1:

Given a hypothesis f , a loss function l , and a sample $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, the *empirical risk* (ER) is defined by

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n l(f(\mathbf{x}_i), y_i),$$

and the *population risk* is defined by

$$L(f) = \mathbb{E}_{(x,y) \sim D}(l(f(x), y)).$$

In supervised learning, we consider the *empirical risk minimization* (ERM) problem:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} L_n(f).$$

If \mathcal{F} is parameterized by θ , i.e., $\mathcal{F} = \{f(\mathbf{x}; \theta) | \theta \in \Lambda\}$, we can also write

$$\hat{\theta}_n := \operatorname{argmin}_{\theta} L_n(\theta),$$

and

$$\boldsymbol{\theta}^* := \operatorname{argmin}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}).$$

We define the *generalization error* to be

$$f(\mathbf{x}; \hat{\boldsymbol{\theta}}_n) - f^* = (f(\mathbf{x}; \hat{\boldsymbol{\theta}}_n) - f(\mathbf{x}; \boldsymbol{\theta}^*)) + (f(\mathbf{x}; \boldsymbol{\theta}^*) - f^*).$$

On the right hand side of the equality, we call the first term *estimation error*, and the second term *approximation error*.

As the ERM problem is sometimes difficult to solve, so we also consider structure ERM:

$$\hat{\boldsymbol{\theta}}_n(\epsilon) = \operatorname{argmin}_{L_n(\boldsymbol{\theta}) \leq \epsilon} \|\boldsymbol{\theta}\|$$

or

$$\hat{\boldsymbol{\theta}}_n(\lambda) = \operatorname{argmin}_{\boldsymbol{\theta}} (L_n(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2).$$

In the following sections, we mainly concentrate on two questions:

- How to choose our hypothesis space?
- how to learn from the chosen hypothesis space?

2 Linear Regression

Let $\mathcal{F} = \{\boldsymbol{\beta}^\top \mathbf{x} | \boldsymbol{\beta} \in \mathbb{R}^d\}$, $l(y, y') = (y - y')^2$, the ERM problem becomes

$$\min \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\beta}^\top \mathbf{x}_i - y_i)^2,$$

let $Z = (\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top)^\top$, the problem can be written as

$$\min \frac{1}{n} \|Z\boldsymbol{\beta} - Y\|^2$$

- For $n < d$, we consider *Ridge regression*:

$$\min \frac{1}{n} \|Z\boldsymbol{\beta} - Y\|^2 + \lambda \|\boldsymbol{\beta}\|^2$$

3 Generalized Linear Model

Let

$$\mathcal{F} = \left\{ \sum_{j=1}^m \beta_j \phi_j(x) \mid \boldsymbol{\beta} \in \mathbb{R}^m \right\},$$

where the ϕ_j 's are basis functions(for example, polynomials, splines, wavelets,...). Write $\Phi = (\phi_j(x_i))_{i,j}$, the ERM problem becomes

$$\hat{\beta} = \operatorname{argmin} \|\Phi\beta - Y\|^2 + \lambda\|\beta\|^2,$$

by differentiating with respect to β , we get

$$\hat{\beta} = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T y.$$

By easy calculations:

$$\hat{\beta} = \Phi^T(\Phi\Phi^T + \lambda I)^{-1}y.$$

Denote $\alpha := (\Phi\Phi^T + \lambda I)^{-1}y$ and consider the case: $m \gg n$

The complexity of calculating the inverse of a m -order matrix is $O(m^3)$ or $O(Cm^2)$. By the above transform, the order of the matrix which will be inverted is decreased to n and the complexity is decreased to $O(n^3)$.

The key observations:

- $\hat{\beta} \in \operatorname{span}\{\phi(x_1), \dots, \phi(x_n)\}$, we solve

$$\min \|\Phi\Phi^T\alpha - Y\|^2 + \lambda\|\Phi^T\alpha\|^2,$$

- The learned model is given by

$$f(x; \hat{\beta}) = \Phi(x)\hat{\beta} = \sum_{i=1}^n \alpha_i \langle \phi(x_i), \phi(x) \rangle$$

Let $k(x, x') = \langle \phi(x), \phi(x') \rangle$, we have

$$f(x; \hat{\beta}) = \sum_{i=1}^n \alpha_i k(x_i, x), \tag{1}$$

with the coefficients are given by

$$\alpha = (K + \lambda I_n)^{-1}Y, \tag{2}$$

where $K = (k(x_i, x_j))_{i,j} \in \mathbb{R}^{n \times n}$. This matrix is called Gram or Kernel matrix.

4 Kernel Method

After the observations in , we have A kernel over \mathcal{X} is a function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ which satisfies two conditions:

- For any $x, x' \in \mathcal{X}$, $k(x, x') = k(x', x)$;
- For any $x_1, x_2, \dots, x_n \in \mathcal{X}$, the Gram matrix $K = (k(x_i, x_j))$ is positive semidefinite.

Kernel Ridge Regression: For any kernel k , the estimator (1) is called kernel ridge regression.

Example:

- $k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}'), k(\mathbf{x}, \mathbf{x}') = k_0(\|\mathbf{x} - \mathbf{x}'\|_2).$

- $k(u, v) = (\langle u, v \rangle + 1)^m$ Denote $u_0 = v_0 = 1$, then

$$k(u, v) = \sum_{J \in \{0,1,\dots,d\}^m} \left(\prod_{i=1}^m u_{J_i} \right) \left(\prod_{i=1}^m v_{J_i} \right) = \langle \phi(u), \phi(v) \rangle$$

- $k(x, x') = e^{-\frac{(x-x')^2}{2}}$ From the Taylor expansion, we have

$$k(x, x') = \sum_{n=0}^{\infty} \frac{1}{\sqrt{n!}} e^{-\frac{x^2}{2}} x^n \frac{1}{\sqrt{n!}} e^{-\frac{x'^2}{2}} x'^n = \sum_{n=0}^{\infty} \phi_n(x) \phi_n(x')$$

5 Random Fourier Feature model

Shift invariant kernel is equivalent to *Random Fourier feature*. So for a high dimensional problem, it's effective to use the latter to decrease the dimension and therefore *Kernel Method* could be used. The argument is promised by the below theorem.

Theorem 5.1 (Bochner). Assume $k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}')$. Then k is a continuous kernel function iff k_0 is the Fourier transform of a non-negative measure.

The probability measure is an example:

$$k_0(\mathbf{x} - \mathbf{y}) = \int_{\mathbb{R}^d} p(\boldsymbol{\omega}) e^{i\boldsymbol{\omega}^T(\mathbf{x}-\mathbf{y})} d\boldsymbol{\omega} = \mathbb{E}_{\boldsymbol{\omega} \sim \mathbb{P}}[z_{\boldsymbol{\omega}}(\mathbf{x}) \bar{z}_{\boldsymbol{\omega}}(\mathbf{y})]. \text{ Where } z_{\boldsymbol{\omega}}(\mathbf{x}) = e^{i\boldsymbol{\omega}^T \mathbf{x}}$$

The vector form is $z_{\boldsymbol{\omega}}(\mathbf{x}) = \begin{pmatrix} \cos(\boldsymbol{\omega}^T \mathbf{x}) \\ \sin(\boldsymbol{\omega}^T \mathbf{x}) \end{pmatrix}$

An example for the theorem: $z_{\boldsymbol{\omega}}(\mathbf{x}) = \sqrt{2} \cos(\boldsymbol{\omega}^T \mathbf{x} + b), b \sim u([0, 2\pi])$

$$\frac{1}{2\pi} \int_0^{2\pi} z_{\boldsymbol{\omega}}(\mathbf{x}, b) z_{\boldsymbol{\omega}}(\mathbf{y}, b) db = \frac{2}{2\pi} \int_0^{2\pi} \cos(\mathbf{w}^T \mathbf{x} + b) \cos(\mathbf{w}^T \mathbf{y} + b) db \quad (3)$$

$$= \frac{1}{\pi} \int_0^{2\pi} \cos(\mathbf{w}^T(\mathbf{x} + \mathbf{y}) + 2b) + \cos(\mathbf{w}^T(\mathbf{x} - \mathbf{y})) db \quad (4)$$

$$= 2 \cos(\mathbf{w}^T(\mathbf{x} - \mathbf{y})) \quad (5)$$

The random approximation of the kernel:

$$k(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\boldsymbol{\omega} \sim \mathbb{P}}[z_{\boldsymbol{\omega}}(\mathbf{x}) z_{\boldsymbol{\omega}}(\mathbf{y})] \simeq \frac{1}{m} \sum_{k=1}^m z_{\boldsymbol{\omega}_k}(\mathbf{x}) z_{\boldsymbol{\omega}_k}(\mathbf{y})$$

The difference between *Emperical Risk* and *Population Risk* could be estimated by the Central limit theorem.

$$\sup_{x,y \in M} |k(\mathbf{x}, \mathbf{y}) - \frac{1}{m} \sum_{k=1}^m z_{\omega_k}(\mathbf{x}) z_{\omega_k}(\mathbf{y})| \sim \frac{C}{\sqrt{m}} = \epsilon, C(\text{diam}(M), d)$$

One could refer to *Random Features for Large-Scale kernel Machines*(2007, Ali Rahimi, Benjamin Recht) for more details.