

Lecture 5 — 13/08

Lecturer: Arjun Bhagoji

Scribe: Aditya Agrawal

5.1 Input Representations

We start off by looking at whether the way we represent data for classification good enough?

A *feature map* $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p$ is simply a map which maps the input data into another space known as the *feature* or *embedding* space.

The k-class classification problem asks for a *map* $f : X \rightarrow [0, 1]^k$ to denote *confidence* of a particular input lying in some class.

The standard classification algorithms directly operate on the input data while *representational learning* first maps raw data into *feature space* on which standard algorithms are applied. Representational learning is generally of two types:

- (a) *Fixed* representations: These are explicit maps which are used on basis on some prior knowledge of data distribution
- (b) *Learned* representations: These are maps which are *learnt* through Neural Networks, or other such techniques

A *kernel* is defined as a dot product in feature space:

$$k : X \times X \rightarrow \mathbb{R} \quad k(x, x') = \langle \phi(x), \phi(x') \rangle$$

The standard dot product arises when we use the linear feature map $\phi(x) = x$.

5.1.1 Feature maps from kernels

We ask the question of when do arbitrary kernel functions arise as a dot product between some feature vectors.

Theorem 5.1. *If $k(x, x')$ is real valued and positive semi definite, then there exists an inner product space \mathcal{X} and feature map $\phi : \mathbb{R}^d \rightarrow \mathcal{X}$ such that $k(x, x') = \langle \phi(x), \phi(x') \rangle$*

We simply define \mathcal{X} as the set of all maps $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with an appropriate inner product and the feature map $\phi : \mathbb{R}^d \rightarrow \mathcal{X}$ as $\phi(x)(\cdot) = k(x, \cdot)$. We also have to validate the existence of the inner product on the function space but the positive semi definiteness of the product will allow us to do so.

5.1.2 Examples of Kernels

- (a) Polynomial Kernels : $k(x, x') = \langle x, x' \rangle^d$
- (b) Gaussian Kernels : $k(x, x') = e^{-\|x - x'\|^2 / 2\sigma^2}$

5.1.3 The Kernel Trick

Well, the representations we dealt with earlier either fixed or learned were into finite dimensional spaces. But fortunately, in all our optimization equations, the mentions of $\phi(x)$ were all as dot products of the form $\phi(x)^T \phi(x')$, which can instead be replaced by $\tilde{k}(x, x')$ leading to a feature map $\tilde{\phi} : \tilde{\phi}(x)(\cdot) = \tilde{k}(x, \cdot)$ where this feature map could be infinite dimensional!!

5.2 Anomaly Detection

We come back to the original question of anomaly detection and propose another method.

5.2.1 One Class SVM

Consider the following constrained parametric optimization problem:

$$\begin{aligned} \min_{w, \gamma, \rho} \quad & \frac{1}{2} \|w\|^2 + \frac{1}{vn} \left(\sum_{i=1}^n \gamma_i \right) - \rho \\ \text{subject to} \quad & \begin{cases} \langle w, \phi(x_i) \rangle \geq \rho - \gamma_i & i = 1, \dots, n, \\ \gamma_i \geq 0 & i = 1, \dots, n, \end{cases} \end{aligned}$$

where v is a hyperparameter which asymptotically represents the fraction of anomalies. At test time, we define the hypothesis maps f, h as the following:

$$\begin{aligned} f_{\theta}(x) &= \langle w, \phi(x) \rangle - \rho \\ h_{\theta}(x) &= \text{sgn}(f_{\theta}(x)) \end{aligned}$$

Note: Finally, the map h is the one which maps the input into a class 0 or 1 where 0 indicates an anomaly.

5.2.2 HW Problems

1. Derive the primal optimization problem of the *lagrangian* of a standard one class SVM machine for both the soft and the hard cases using the KKT conditions
2. Derive the same for the v -SVM as well
3. Derive the dual of these primal objective functions as well

We provide the solution to the one of the *lagrangian* objective function mentioned above, the rest of them are left as exercises to the student, by using the KKT conditions.

We are given:

$$\begin{aligned} \langle w, \phi(x_i) \rangle &\geq \rho - \gamma_i & i = 1, \dots, n, \\ \gamma_i &\geq 0 & i = 1, \dots, n, \end{aligned}$$

We use the set of variables α_i to model $\rho - \gamma_i - w^T \phi(x_i) \leq 0$, the set η_i to model $-\gamma_i \leq 0$ and μ to model $-\rho \leq 0$. The lagrangian we then get is:

$$\begin{aligned} \mathcal{L}(w, b, \rho, \gamma, \alpha, \eta, \mu) = & \frac{1}{2} \|w\|^2 - \rho + \frac{1}{vn} \sum_{i=1}^n \gamma_i + \sum_{i=1}^n \alpha_i \left(\rho - \gamma_i - y_i (w^T \phi(x_i) + b) \right) - \sum_{i=1}^n \eta_i \gamma_i - \mu \rho \\ & \alpha_i \geq 0, \eta_i \geq 0, \mu \geq 0 \end{aligned}$$

On applying the KKT conditions to these, we get the dual objective function as:

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq \frac{1}{vn}, \quad i = 1, \dots, n, \\ & \sum_{i=1}^n \alpha_i y_i = 0, \\ & \sum_{i=1}^n \alpha_i \geq 1 \quad (\text{and } \sum_{i=1}^n \alpha_i = 1 \text{ when } \rho > 0). \end{aligned}$$

5.2.3 SLT- like Guarantees on the One Class SVM

Theorem 5.2. (informally) Let $R_{w,\gamma} = \{x | f_{\theta}(x) \geq 0\}$. Then for all $\delta > 0$, with probability $1 - \delta$ over the draws $\{x_i\}_{i=1}^n$, for any $B > 0$, we get

$$\mathbb{P}_{x \sim \mathbb{P}^+} [x | x \notin R_{w,\rho-B}] \leq \frac{2}{n} \left(k(n, B, \|W\|, \rho) + \log_2 \left(\frac{n^2}{2\delta} \right) \right)$$

where k is some non-negative function of problem parameters

5.2.4 A Unifying View of Anomaly Detection

In essence, anomaly detection involves designing *score functions* which work well during test time. When the score of a particular sample is high enough, we deem it to be from our sample distribution. Students can check out A Unified Survey on Anomaly, Novelty, Open-Set, and Out-of-Distribution Detection: Solutions and Future Challenges for further reading.

5.3 OOD Detection

The OOD problem or the *Out Of Distribution* problem is the supervised case of the anomaly detection problem where we are now also provided labels on the input draws. Our key question is to understand how OOD detection is done nowadays.

5.3.1 Deep OCSVM or Deep SVDD

In this methodology, we use a learned feature map $f_\theta(x)$, where θ parameterizes some neural network, with θ learned jointly over the supervised learning component and the anomaly detection (R, c) or (w, γ, ρ) . The anomaly detection optimization problem is given as

$$\min_{R, c, \theta} R^2 + \frac{1}{\nu n} \sum_{i=1}^n \max(0, \|\phi_\theta(x_i) - c\|^2 - R^2) + \text{Reg}(\theta)$$

Here, typically some regularization is put on θ since otherwise it tends to fit maps such that all points in feature space start falling within the sphere. It is typically trained through stochastic gradient descent

5.3.2 MSP and ODIN

MSP stands for *Maximum Softmax Probability* and works on the principle that the NN holds enough discriminative power to classify in distribution samples confidently. This leads to the formulation where we call an input out of distribution when the NN cannot classify it confidently.

ODIN or *Out of Distribution for Neural Networks* enhances MSP by preprocessing data smartly and adding a temperature control on the softmax performed. It works in the following way:

1. Preprocesses input x as

$$x' = x + \varepsilon * \text{sgn} \left(\nabla_x \log(f_\theta^{\hat{y}}(x, \tau)) \right)$$

$$f_\theta^i(x, \tau) = \frac{e^{\frac{l_i(x)}{\tau}}}{\sum_{i=1}^n e^{\frac{l_i(x)}{\tau}}}$$

where $l_i(x)$ is the logit output of neural network $l : \mathbb{R}^d \rightarrow \mathbb{R}^p$

2. Defines the anomaly hypothesis map h as

$$s_\theta(x) = \max_i f_\theta^i(x, \tau)$$

$$h_\theta(x) = 2 * (\mathbf{1}[s_\theta(x) < \delta] - 1) - 1$$

5.3.3 Outlier Exposure

Well the issue with our current methods is that they are overfitted on in distribution data. It is very likely that out of distribution data looks very different and might lead the overfit classifier to falsely label them as being in distribution. To prevent this, we perform outlier exposure by slightly tweaking our NN loss functions. Our goal is to now minimize

$$\min_{\theta} \mathbb{E}_{x \sim \mathbb{P}^+} l(h_\theta(x), y) + \lambda \mathbb{E}_{x \sim \mathbb{P}^-} l_{OE}(f_\theta(x))$$

Here λ controls the rate of outliers we expect in testing and l_{OE} is the loss we want to calculate over outlier data. We generally assume that the out of distribution is uniform in some sense.

5.3.4 Enhancing MSP

We essentially want to temper the confidence of our classifier, by saying "Hey! Don't be so confident about your predictions, they might be on out of distribution data!!". We design two outlier two loss functions for the same.

1.

$$l_{OE}(x) = - \sum_{j=1}^{|C|} \frac{1}{|C|} \log(f_{\theta}^j(x))$$

We *assume* that a good classification is a confident classification, in the sense that if our classifier classified something well, then we inherently assumed such a point belonged to our distribution

2. To make a distinction between a *good* and *confident* classification we artificially introduce a confidence estimation branch $c(x) \in [0, 1]$ and set $l_{OE}(x) = -\log(c(x))$.

During training we then try to minimize

$$l(x) = - \sum_{j=1}^{|C|} \log(cf_{\theta}(x_j) + (1-c)y_j)y_j - B\log(c(x))$$

Essentially, what we try to say is that if we are not confident about whether the input is from distribution, then our NN doesn't pay a training loss, however we pay a confidence loss. If we are confident, then we incur a training loss as well. This flips our notion from good \implies confident to good \Leftarrow confident

5.3.5 Grad Norm KL Divergence

Instead of looking at the maximum softmax probability as a measure of confidence, we look at the gradient of the norm of the KL divergence as such a measure. The intuition is that low gradient norm implies points close to margins and so more likely to be OOD. We define this using

$$\begin{aligned} D_{KL}(U||f_{\theta}(x)) &= \sum_{j=1}^{|C|} \frac{1}{|C|} \log\left(\frac{1/|C|}{f_{\theta}^j(x)}\right) \\ &= -\frac{1}{|C|} \sum_{j=1}^{|C|} \log(|C|f_{\theta}^j(x)) \end{aligned}$$

Now,

$$\nabla_{\theta} D_{KL}(U||f_{\theta}(x)) = -\frac{1}{|C|} \sum_{j=1}^{|C|} \frac{\nabla_{\theta} f_{\theta}(x_i)}{f_{\theta}^j(x)}$$

If the point is far away from the margins, then some value of $f_{\theta}(x_i)$, will be small, leading to high gradient norm. Conversely, if the point is close to margin, then the gradient norm will be smaller.

Our hypothesis function h then is $h_{\theta}(x) = 2 * (\mathbf{1}[\|\nabla_{\theta} D_{KL}\|_2 < \delta] - 1) - 1$

5.3.6 Cons of the above methods

While the above methods are very intuitive and make sense, most of their guarantees are empirical and not theoretical. They have also been designed in an adhoc manner with not much theoretical basis as to why they might be good. It might be useful to try and get some SLT like bounds on these like the one we had for the one class SVM.