

Lecture 3 — 09/13

Lecturer: Dimitris Papailiopoulos

Scribe: Ananth Sridhar, Ashwin Varadarajan

Note: These lecture notes are still rough, and have only have been mildly proofread.

3.1 Review

- The Empirical Risk concentrates for **finite hypothesis classes**

– *Floating Point Assumption:*

$$n_{\text{samples}} = O\left(\frac{n_{\text{parameters}}}{\epsilon^2}\right) \quad (3.1)$$

- Empirical risk and true risk (concentration). If, h is a predictor,

$$\hat{R}[h] \leq R[h] + \epsilon \quad (3.2)$$

- Earlier, we assumed that the cardinality of the hypothesis class H is finite. Can we handle infinite hypothesis classes? Yes, using the VC-dimension (section 3.2).

3.2 VC dimension

- From Vapnik-Chervonenkis Theory [1971]
- The VC-dimension expresses the predictive power of a classifier
- In the general case, finding the VC-dimension of a hypothesis class is not trivial, and we only know the VC-dimensions for a handful of hypothesis classes. However, there are a few special cases of hypothesis classes where the VC-dimension is easy to calculate (see example). Also, we can always bound the VC-dimension if the Hypothesis class is a finite set, as in (eqn 3.3).

If \mathcal{H} is the hypothesis class,

$$\text{VC}_{\text{dimension}}[\mathcal{H}] \leq \log |\mathcal{H}| \quad (3.3)$$

The above inequality implies that the VC dimension is a stronger notion of complexity than simply the number of predictors.

Example: if h is a predictor, and \mathcal{H} is the hypothesis class

$$\text{If, } h(x_i) = y_i \quad \forall \quad i \in \{1 \dots n\}, \quad \text{for any } n \\ \text{VC}_{\text{dimension}}[\mathcal{H}] = \infty$$

$$\text{If, } \mathcal{H} = \{\text{All Hyperplanes in } d \text{ dimensions}\} \\ \text{VC}_{\text{dimension}}[\mathcal{H}] = d + 1$$

3.3 Concentration of Empirical Risk

The following holds with probability $1 - \delta$, for finite and infinite hypothesis classes.

$$\sup_{h \in \mathcal{H}} \left| \hat{R}_n[h] - R[h] \right| \leq O \left(\sqrt{\frac{\text{VC}_{\text{dimension}}[\mathcal{H}] \log \left(\frac{n}{\text{VC}_{\text{dimension}}[\mathcal{H}]} \right) + \log \frac{1}{\delta}}{n}} \right) \quad (3.4)$$

3.3.1 Concentration of Empirical Risk for ERM

$$\hat{h}^* \in \arg \min_{h \in \mathcal{H}} \hat{R}_n[h] \quad (\text{Empirical Risk Minimizer}) \quad (3.5)$$

$$h^* \in \arg \min_{h \in \mathcal{H}} R[h] \quad (\text{True Risk Minimizer}) \quad (3.6)$$

$$R[\hat{h}^*] - R[h^*] = ? \quad (\text{How close are we?}) \quad (3.7)$$

For the purposes of the lectures, we will assume that the set containment is an equality, that is, minimizers are unique.

Theorem 3.1. *The empirical risk of the empirical risk minimizer concentrates around the true risk of the true risk minimizer*

Proof:

$$R[\hat{h}^*] = \begin{cases} R[\hat{h}^*] + (R[h^*] - R[h^*]) \\ \quad + (\hat{R}_n[h^*] - \hat{R}_n[h^*]) \\ \quad + (\hat{R}_n[\hat{h}^*] - \hat{R}_n[\hat{h}^*]) \end{cases} \quad (3.8)$$

$$R[\hat{h}^*] - R[h^*] = \begin{cases} (\hat{R}_n[h^*] - R[h^*]) & \text{concentrates: } \leq \epsilon \\ + (R[\hat{h}^*] - \hat{R}_n[\hat{h}^*]) & \text{concentrates: } \leq \epsilon \\ + (\hat{R}_n[\hat{h}^*] - \hat{R}_n[h^*]) & \text{concentrates: } \leq 0 \end{cases} \quad (3.9)$$

$$R[\hat{h}^*] - R[h^*] \leq 2\epsilon \quad (3.10)$$

□

3.4 Examples of learning problems

The following are some typical classes of problems in Machine Learning:

3.4.1 Linear Regression

Fitting a line or hyperplane to data points so as to minimize the sum of distance of predicted labels to the observed ones. Involves minimizing the norm:

$$\min \frac{1}{n} \sum_{i=1}^n (x_i^T w - y_i)^2 \quad (3.11)$$

3.4.2 Binary Classification

Assigning binary labels to data, hence the variable y is discrete and $y|x$ is a Bernoulli distribution. Involves solving the logistic regression:

$$\min \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) \quad (3.12)$$

3.4.3 Support Vector Machines

Involves minimizing the maximum margin from a classifier, with a regularization term.

$$\min \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i x_i^T w) + \lambda \|w\| \quad (3.13)$$

3.4.4 Feed Forward Neural Network

Each layer involves a linear sum of weighted features, passed through an activation function σ , which is usually a step function or a smooth version of it, such as $\sigma(x) = -1/(1+\exp(-x))$. This can be repeated over a number of hidden layers between the input and the output layers (deep networks).

3.5 Useful properties for loss functions

In general, the problem $\min_x f(x)$ is hard to solve as is, and hence we use additional properties of the loss function f to speed up computation.

3.5.1 Convex Functions

Any local minimum of the function is also a global minimum.

$$f(\alpha\vec{x} + (1 - \alpha)\vec{y}) \leq \alpha f(\vec{x}) + (1 - \alpha)f(\vec{y}) \quad (3.14)$$

$$\forall \alpha \in [0, 1], \quad \vec{x}, \vec{y}$$

3.5.2 L-Lipschitz functions

The function does not change too fast.

$$|f(\vec{x}) - f(\vec{y})| \leq L\|\vec{x} - \vec{y}\| \quad (3.15)$$

3.5.3 β -smooth functions

The function has a Lipschitz gradient.

$$\|\nabla f(\vec{x}) - \nabla f(\vec{y})\| \leq \beta\|\vec{x} - \vec{y}\| \quad (3.16)$$

3.5.4 λ -strongly convex functions

The function has a unique global minimum.

$$f(a\vec{x} + (1 - a)\vec{y}) \leq af(\vec{x}) + (1 - a)f(\vec{y}) - \frac{\lambda}{2}a(1 - a)\|\vec{x} - \vec{y}\|^2 \quad (3.17)$$

$$\forall a \in [0, 1], \quad \vec{x}, \vec{y}$$

In other words, f is λ -strongly convex **iff**

$$f(\vec{x}) - \lambda\|\vec{x}\|^2 \quad \text{is convex} \quad (3.18)$$

This property is especially nice, since we can achieve this for ERM by regularization. Strong convexity gives improved algorithmic convergence rates (the algorithms make use of the regularization for faster convergence). Note: All convex functions are 0-strongly convex (trivial case).

3.6 Additional Reading

- A nice reference on the VC-dimension (explained with an example): [Link](#)
- An additional resource for insights into strong convexity and some settings where it is utilized: [Link](#)