

Lecture 12 - 21-04-2020

1.1 Non parametric algorithms

We talk about **consistency**: as the training size grows unbounded the expected risk of algorithms converge to Bayes Risk.

Now we talk about **non parametric algorithm**: the structure of the model is determined by the data.

Structure of the model is fixed, like the structure of a Neural Network but in non parametric algorithm will change structure of the model as the data grows (K_{NN} and tree predictor).

If I live the tree grow unboundedly then we get a non parametric tree, but if we bound the grows then we get a parametric one.

The converge rate of Bayes Risk (in this case doubled) was small. Converge of 1-NN to $2\ell_D(f^*)$ is $m^{-\frac{1}{d+1}}$ so we need an exponential in the dimension. And we need this is under Lips assumption of η .

It's possible to converge to Bayes Risk and it's called **No free lunch**.

1.1.1 Theorem: No free lunch

Let a sequence of number $a_1, a_2 \dots \in \mathbb{R}$ such that they converge to 0.

Also $\frac{1}{22222222} \geq a_1 \geq a_2 \geq \dots \forall A$ for binary classification $\exists D$ s. t.

$\ell_D(f^*) = 0$ (zero-one loss) so Bayes risk is zero and $\mathbb{E}[\ell_D(A(S_M))] \geq a_m \quad \forall m \geq 1$

Any Bayes Optimal you should be prepared to do so on long period of time. This means that:

- For specific data distribution D , then A may converge fast to Bayes Risk.
- If η is Lipschitz then it is continuous. This mean that we perturb the input by the output doesn't change too much.
- If Bayes Risk is 0 ($\ell_D(f^*) = 0$) function will be discontinuous

This result typically people think twice for using consistent algorithm because

I have Bayes risk and some non consistent algorithm that will converge to some value ($\ell_D(\hat{h}^*)$). Maybe i have Bayes risk and the convergence takes a

1.2 Highly Parametric Learning Algorithm

1.2.1 Linear Predictors

Our domain is Euclidean space (so we have points of numbers).

$$X \text{ is } \mathbb{R}^d \quad x = (x_1, \dots, x_d)$$

A linear predictor will be a linear function of the data points.

$$h : \mathbb{R}^d \longrightarrow Y \quad h(x) = f(w^T x) \quad w \in \mathbb{R}^d$$

$$f : \mathbb{R} \longrightarrow Y$$

And this is the dot product that is

$$w^T x = \sum_{i=1}^d w_i x_i = \|w\| \|x\| \cos \Theta$$

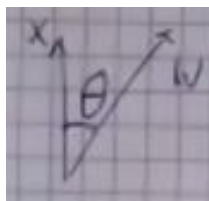


Figure 1.2: Dot product

Suppose we look a regression with square loss.

$$Y = \mathbb{R} \quad h(x) = w^T x \quad w \in \mathbb{R}^d$$

$$f^*(x) = \mathbb{E}[Y|X = x]$$

Binary classification with zero-one loss $Y = \{-1, 1\}$ We cannot use this since is not a real number but i can do:

$$h(x) = \text{sgn}(w^T x) \quad \text{sgn}(x) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0 \end{cases}$$

where sgn is a sign function. Linear classifier.

$\|X\| \cos \Theta$ is the length of the projection of x onto w

Now let's look at this set:

$$\{x \in \mathbb{R}^d : w^T x = c\}$$

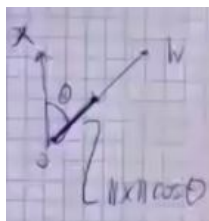


Figure 1.3: Dot product

This is a hyperplane.

$$\|w\| \|x\| \cos \Theta = c \quad \|x\| \cos \Theta = \frac{c}{\|w\|}$$

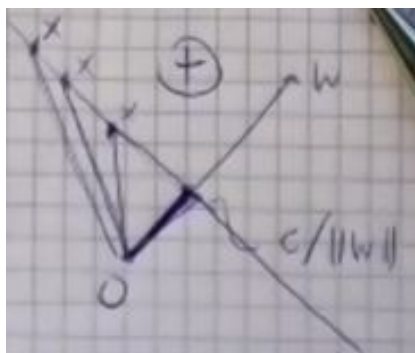


Figure 1.4: Hyperplane

So (w, c) describe an hyperplane.

We can do binary classification using the hyperplane. Any points that lives in the positive half space and the negative. So the hyperplane is splitting in halves. $H \equiv \{x \in \mathbb{R}^d : w^T x = c\}$

$$H^+ \equiv \{x \in \mathbb{R}^d : w^T x > c\} \quad \text{positive } h_s$$

$$H^- \equiv \{x \in \mathbb{R}^d : w^T x \leq c\} \quad \text{negative } h_s$$

$$h(x) = \begin{cases} +1 & \text{if } x \in H^+ \\ -1 & \text{if } x \notin H^+ \end{cases} \quad h(x) = \text{sgn}(w^T x - c)$$

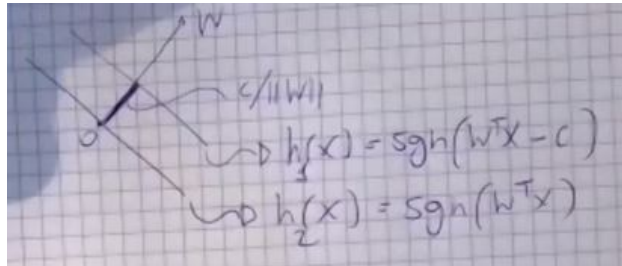


Figure 1.5: Hyperplane

h_1 is non-homogenous linear classifier.
 h_2 is homogenous linear classifier. Any homogenous classifier is equivalent to

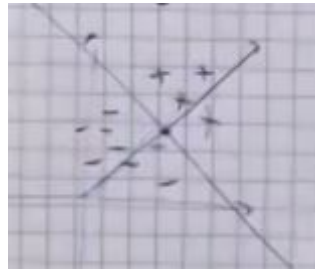


Figure 1.6: Hyperplane

this:

$$\{x \in \mathbb{R}^d : X = c\} \text{ is equivalent to } \{x : \mathbb{R}^{d+1} : \nu^T x = 0\}$$

$$\nu = (w_1, \dots, w_d, -c) \quad x' = (x_1, \dots, x_d, 1)$$

So we added a dimension.

$$w^T x = c \Leftrightarrow \nu^T x' = 0$$

$$\sum_i w_i x_i = c \Leftrightarrow \sum_i w_i x_i - c = 0$$

Rule:

When you learn predictor just add an extra feature to your data points, set it to 1 and forget about non-homogenous stuff.

One dimensional example

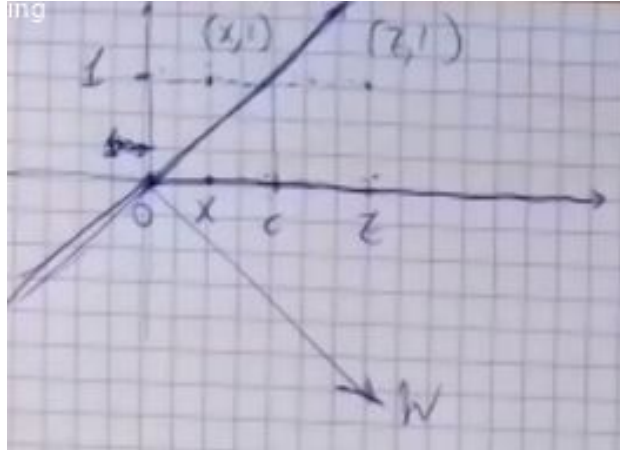


Figure 1.7: Example of one dimensional hyperplane

I have negative (left of $(x, 1)$) and positive point (left of $(z, 1)$) classified

Now i want to learn linear classifier. How can i do it?

$$H_d = \{ h : \exists w \in \mathbb{R}^d h(x) = \text{sgn}(w^T x) \}$$

Parametric!

We expect high bias a low variance.

$$\begin{aligned} \text{ERM} \quad \hat{h}_S &= \arg \min_{h \in H_d} \frac{1}{m} \cdot \sum_{t=1}^m I\{h(x_t) \neq y_t\} = \\ &= \arg \min_{w \in \mathbb{R}^d} \frac{1}{m} \cdot \sum_{t=1}^m I\{y_t w^T x_t \leq 0\} \end{aligned}$$

A bad optimisation problem!

FACT:

It is unlikely to find an algorithm that solves ERM for H_d and zero-one loss efficiently.

NP completeness problems!

It's very unlikely to solve this problem.

This problem is called **MinDisagreement**

1.2.2 MinDisagreement

Instance: $(x_1, y_1) \dots (x_m, y_m) \in \{0, 1\}^d \times \{-1, 1\}$, $k \in \mathbb{N}$

Question: Is there $w \in \mathbb{R}^d$

s.t. $y_t w^T x_t \leq 0$ for at most k indices $t \in \{1, \dots, m\}$

This is NP-complete!