

Introduction to ML

September 29th, 2021

Model Selection & Generalization

- Learning is an **ill-posed problem**; data is not sufficient to find a unique solution
 - Every sample is boolean (d-dimension 2^d possible input function)
 - Every output is Boolean (input-output pair of hypotheses function left 2^{2^d})
 - Each data sample we kills half (the rest are consistent – different function gives the same answer)
 - With N samples there remains $2^{2^d - N}$
 - Impossible for a UNIQUE answer, we usually only see a fraction of data samples

All possible hypothesis functions that takes tuple binary input to binary output

- we have a 2^d - binary input (0,1), output is 0 or 1
- there are 2^d possible combination of binary **input** (this is a set with 2^d elements)
- we have another set that is has 2 element (yes/no)
- How many ways we can take 2^d elements to 2 possible outcomes?
 - should be 2^{2^d}

If you have N samples observed, how many possible options are left:

- $2^{2^d - N}$

- take $d=2$ as an example:
- first possible function (left input, right output)
 - (0,0)|0
 - (0,1)|0
 - (1,0)|0
 - (1,1)|0
- second possible function
 - (0,0)|1
 - (0,1)|0
 - (1,0)|0
 - (1,1)|0
- third possible function
 - (0,0)|0
 - (0,1)|1
 - (1,0)|0
 - (1,1)|0
- fourth possible function
 - (0,0)|0
 - (0,1)|0
 - (1,0)|1
 - (1,1)|0
- ... so on, there are total of $2^4 \rightarrow 2^{2^2} \rightarrow 16$ functions

Process of of a Supervised Learner

1. Model:

$$g(x|\theta)$$

2. Loss function:

$$E(\theta|\mathcal{X}) = \sum_t L(r^t, g(x^t|\theta))$$

3. Optimization procedure:

$$\theta^* = \arg \min_{\theta} E(\theta|\mathcal{X})$$

CHAPTER 3:

Bayesian Decision Theory

Classification

Observation \rightarrow measurable input

- Credit scoring: Inputs are income and savings.

Output is low-risk vs high-risk

- Input: $\mathbf{x} = [x_1, x_2]^T$, Output: $C \in \{0, 1\}$

Bernoulli rv
Condition on
 x_1, x_2

- Prediction:

choose $\begin{cases} C = 1 \text{ if } P(C = 1 | x_1, x_2) > 0.5 \\ C = 0 \text{ otherwise} \end{cases}$

or

choose $\begin{cases} C = 1 \text{ if } P(C = 1 | x_1, x_2) > P(C = 0 | x_1, x_2) \\ C = 0 \text{ otherwise} \end{cases}$

ERROR? $1 - \max (P(C=1|x_1, x_2), P(C=0|x_1, x_2))$

Maximum A-posterior decision rule (find the class with highest probability)

Guarantee the error is smallest as measured in terms of probability 6

Bayes' Rule

What is C (high risk/low risk) given the two inputs

posterior

prior

likelihood

$$P(C | \mathbf{x}) = \frac{P(C) p(\mathbf{x} | C)}{p(\mathbf{x})}$$

evidence

$$P(C = 0) + P(C = 1) = 1$$

$$p(\mathbf{x}) = p(\mathbf{x} | C = 1)P(C = 1) + p(\mathbf{x} | C = 0)P(C = 0)$$

$$p(C = 0 | \mathbf{x}) + p(C = 1 | \mathbf{x}) = 1$$

Total
probability
theorem

Bayes' Rule: $K > 2$ Classes

$$\begin{aligned} P(C_i | \mathbf{x}) &= \frac{p(\mathbf{x} | C_i)P(C_i)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k)P(C_k)} \end{aligned}$$

$$P(C_i) \geq 0 \text{ and } \sum_{i=1}^K P(C_i) = 1$$

choose C_i if $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

Make a
'classification'
decision by looking
at the posterior
distribution

Losses and Risks

- Actions: α_i
 - This concept is important in cases when ‘loss’ associated to each action may not be the same
 - Finance, medical, ...so on
- Loss of α_i when the state is $C_k : \lambda_{ik}$
- **Expected risk** (Duda and Hart, 1973)

Loss of making a decision to assign input to class i, when the true class is k

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x})$$

choose α_i if $R(\alpha_i | \mathbf{x}) = \min_k R(\alpha_k | \mathbf{x})$

Minimizing expected risk by picking that action i

Losses and Risks: 0/1 Loss

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

Correct
decision: no
loss
Incorrect
decision : 1
loss

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x})$$

$$= \sum_{k \neq i} P(C_k | \mathbf{x})$$

All class except
that correct one

$$= 1 - P(C_i | \mathbf{x})$$

For minimum risk, choose the most probable class

Losses and Risks: Reject

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = K + 1, \quad 0 < \lambda < 1 \\ 1 & \text{otherwise} \end{cases}$$

Additional action when in doubt

Loss incurred when in doubt
(reject as another action)

$$R(\alpha_{K+1} | \mathbf{x}) = \sum_{k=1}^K \lambda P(C_k | \mathbf{x}) = \lambda$$

Probability sums to 1

$$R(\alpha_i | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$$

Risk needs to be minimized

choose C_i if $P(C_i | \mathbf{x}) > P(C_k | \mathbf{x}) \quad \forall k \neq i$ and $P(C_i | \mathbf{x}) > 1 - \lambda$
reject otherwise

Has to be greater than reject

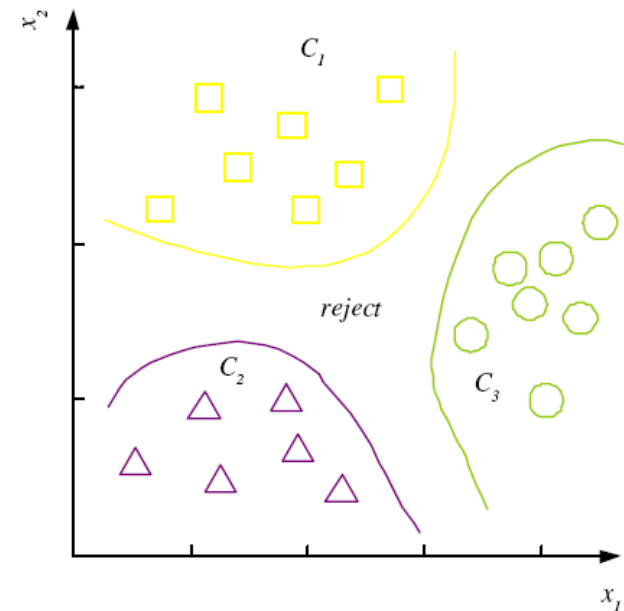
Discriminant Functions

Classification rule: pick one such function that maximizes

choose C_i if $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$

$$g_i(\mathbf{x}) = \begin{cases} -R(\alpha_i | \mathbf{x}) \\ P(C_i | \mathbf{x}) \\ p(\mathbf{x} | C_i)P(C_i) \end{cases}$$

Three different discriminant functions



K decision regions $\mathcal{R}_1, \dots, \mathcal{R}_K$

$$\mathcal{R}_i = \{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$$

Divides feature space into K region
For those inputs \mathbf{x} , find the function that give the largest value (use that function to carve out a region)


$K=2$ Classes

- Dichotomizer ($K=2$) vs Polychotomizer ($K>2$)
- $g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$

$$\text{choose } \begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

- *Log odds:*

Needs only 1 value to
make a decision



$$\log \frac{P(C_1 | \mathbf{x})}{P(C_2 | \mathbf{x})}$$

A discriminant function



Utility Theory


- Prob of state k given evidence \mathbf{x} : $P(S_k | \mathbf{x})$
- Utility of α_i when state is k : U_{ik}
- Expected utility:

$$EU(\alpha_i | \mathbf{x}) = \sum_k U_{ik} P(S_k | \mathbf{x})$$

Choose α_i if $EU(\alpha_i | \mathbf{x}) = \max_j EU(\alpha_j | \mathbf{x})$

Loss: negative utility

Pretty much terminology
and used quite often in
game theory, economics
so on



Association Rules

- Association rule: $X \rightarrow Y$
- *People who buy/click/visit/enjoy X are also likely to buy/click/visit/enjoy Y.*
 - Dependency, it's not cause and effect but good enough for making a business decision
- A rule implies association, not necessarily causation.

Association measures

Joint probability,
Make this large (increase the basis)
Significance of the rule
Useless if this number is low

- Support ($X \rightarrow Y$):

$$P(X, Y) = \frac{\#\{\text{customers who bought } X \text{ and } Y\}}{\#\{\text{customers}\}}$$

- Confidence ($X \rightarrow Y$):

$$P(Y | X) = \frac{P(X, Y)}{P(X)}$$

- conditional probability
- Should be larger than $P(Y)$
- **Strength** of the association rule

- Lift ($X \rightarrow Y$):

$$= \frac{\#\{\text{customers who bought } X \text{ and } Y\}}{\#\{\text{customers who bought } X\}}$$

$$= \frac{P(X, Y)}{P(X)P(Y)} = \frac{P(Y | X)}{P(Y)}$$

- If independent, this is 1
- If lift > 1 X makes Y more likely

Example

Transaction	Items in basket
1	milk, bananas, chocolate
2	milk, chocolate
3	milk, bananas
4	chocolate
5	chocolate
6	milk, chocolate

SOLUTION:

milk \rightarrow bananas : Support = 2/6, Confidence = 2/4

bananas \rightarrow milk : Support = 2/6, Confidence = 2/2

milk \rightarrow chocolate : Support = 3/6, Confidence = 3/4

chocolate \rightarrow milk : Support = 3/6, Confidence = 3/5

In making a decision, take ‘high support + high confidence rule’