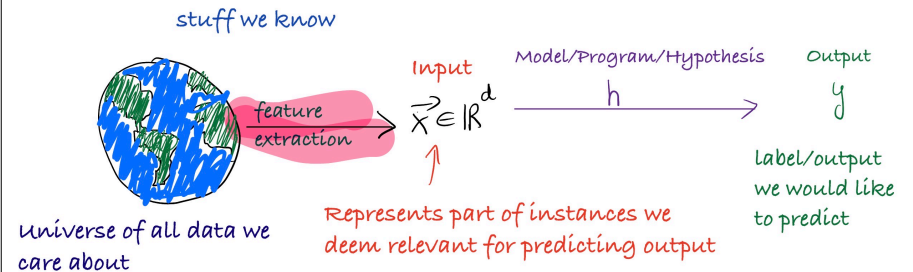


CS 3/5780: Introduction to Machine Learning

Lecture 2: Supervised Learning Setup

SUPERVISED LEARNING



FEATURE VECTORS

What feature vectors would you use for the following tasks:

- ① Classify if a given email is spam or not
- ② Predict the microsoft stock price for the next day
- ③ Where will jupiter be in the night sky tomorrow?
- ④ Is there a pedestrian in front of the car?

FINDING A GOOD HYPOTHESIS

- Traditional approach: Pay a programmer to program h
- Machine Learning: Learn h from examples or Data

SUPERVISED LEARNING SETUP

- Data $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- Loss function $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$ measuring cost of error
- Model class \mathcal{H} consisting of hypotheses or models $h : \mathcal{X} \mapsto \mathcal{Y}$
- Learning algorithm that picks hypothesis $h \in \mathcal{H}$ based on data D

DATA

- $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ where x_i is the d -dimensional input feature vector for the i 'th example and $y_i \in \mathcal{Y}$ is the output for i 'th example
- Scenarios for \mathcal{Y}
 - Binary classification $\mathcal{Y} = \{0, 1\}$
 - Multiclass Classification $\mathcal{Y} = \{1, 2, \dots, K\}$
 - Regression \mathcal{Y} is the set of real numbers
 - Ranking \mathcal{Y} is ordered set of choices
 - Structured prediction \mathcal{Y} Eg. protein structure
- Feature vectors
 - Text: Bag of words features
 - Images: Raw pixels flattened to a vector

LOSS FUNCTION

- Classification loss $\ell(y, y') = \mathbf{1}\{y' \neq y\}$
- Squared loss $\ell(y', y) = (y' - y)^2$
- Absolute loss $\ell(y', y) = |y' - y|$

GENERALIZATION

- Idea: Find a model with low loss on Data D
- Does this always work?
- Eg. Algorithm Memorizer

TEST/TRAIN SPLIT

- Split Data into training set D_{trn} , test set D_{tst} and validation set D_{val}
- Choose h based off of D_{trn} (+ D_{val})
- Evaluate the algorithm on D_{tst}
- How should we split data?

RISK OR POPULATION LOSS

- We often model data at deployment time as being drawn from a distribution P over input output pair
- iid draw from P often makes sense for non-temporal data
- Risk of a hypothesis:

$$\mathcal{E}(h) = \mathbb{E}_{(x,y) \sim P} [\ell(h(x), y)]$$

- Goal of ML algorithm is to pick hypothesis h based on data D to minimize risk $\mathcal{E}(h)$

PURPOSE OF SPLIT

- What is a good proxy for Risk?
- We use D_{trn} to learn/pick hypothesis h , and use D_{tst} to evaluate this hypothesis
- Why do we need D_{val} ?

NO FREE LUNCH!

- You must make assumptions in order to learn
- Corollary: No algorithm works in all settings
- For every method you encounter, question to ask: what is the underlying assumption needed to make this algorithm work?

Question 1: Spam filter you are given spam mail and labels for the past year

How would you split to test and train?

Question 2: Given hypothesis h learnt using D_{trn} I claim that

$$\mathcal{E}_{D_{\text{ts}}}(h) = \frac{1}{|D_{\text{tst}}|} \sum_{(x,y) \in D_{\text{tst}}} l(h(x), y)$$

Is an unbiased and good estimate (when size of training set is large) for risk, why?