

5: Machine Learning Basics

Saturday, June 9, 2018 5:55 PM

hyperparameters: determined outside the learning algorithm

ML is form of applied statistics:

↑ emphasis on using computation to estimate functions

↓ emphasis on proving confidence intervals

most deep learning algorithms based on stochastic gradient descent
optimization algorithm

(S.1) Learning Algorithms

(S.1.1) The Task: T (what we want to learn to do)

- how the ML system should process an example

example: collection of features - quantitatively measured

$x \in \mathbb{R}^n$ feature vector

Examples of ML tasks:

- Classification: learn $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$

$y = f(x)$ feature vector k output classes
e.g., object recognition

variation: output probability distribution

- Classification with missing inputs:

learn a set of functions, one for each possible set of input features

e.g. - medical diagnosis: only have results of some tests

- Regression: learn $f: \mathbb{R}^n \rightarrow \mathbb{R}$

e.g. - predict expected claim amount

- Transcription: unstructured \rightarrow structured data
(text)

e.g., OCR (image \rightarrow text), speech recognition (sound \rightarrow text)

- Machine translation e.g., English \rightarrow Russian

- Structured output e.g., Parsing, image/video segmentation

- Anomaly detection (really, a kind of classification problem)

- Synthesis and sampling e.g., speech synthesis

- Imputation of missing values

- Denoising corrupted example \rightarrow clear example

- Density estimation or probability mass function estimation

...

(S.1.2) Performance measure: P

classification: accuracy, error rate

for continuous tasks,

average log-probability

expected 0-1 loss

incorrect
correctly classified

often difficult to choose P .

(5.1.3) Experience: E

unsupervised: experience of dataset, but no given labels

attempt to learn structure, probability distribution

Supervised: each example has an associated label or target

difference is not fully defined

$$p(x) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1}) \quad \begin{array}{l} \text{decompose feature vector} \\ \text{turns unsupervised problem of modeling} \end{array}$$

else: supervised learning problem: learn $p(y|x)$

$$p(y|x) = \frac{p(x,y)}{\sum_y p(x,y)} \quad \begin{array}{l} \text{labels feature vector} \\ \text{feature vector} \end{array}$$

reinforcement learning: all agents interact with environment \rightarrow feedback loop
instead of fixed dataset

dataset: design matrix - each row is one example, feature vector (+ label if supervised)
assumes same size for each example

(5.1.4) Example: Linear Regression

Input: $x \in \mathbb{R}^n$

Output: $y \in \mathbb{R}$

scalar

$$\hat{y} = w^T x.$$

prediction

$w \in \mathbb{R}^n$ is learned vector of parameters

w_i is coefficient for x_i

(weights)

measure of performance: mean squared error

$$MSE_{\text{test}} = \frac{1}{m} \sum_i (\hat{y}^{(\text{test})} - y^{(\text{test})})^2 = \frac{1}{m} \left\| \hat{y}^{(\text{test})} - y^{(\text{test})} \right\|_2^2$$

prediction actual

simple ML training algorithm: minimize MSE on training data

$$\nabla_w MSE_{\text{train}} = 0 \quad \text{solve for gradient} = 0$$

$$\Rightarrow \nabla_w \frac{1}{m} \| X^{(\text{train})} w - y^{(\text{train})} \|_2^2 = 0.$$

$$\nabla_w (X^{(\text{train})} w - y^{(\text{train})})^T (X^{(\text{train})} w - y^{(\text{train})}) = 0$$

$$\Rightarrow \nabla_w (w^T X^{(\text{train}) T} X^{(\text{train})} w - 2w^T X^{(\text{train}) T} y^{(\text{train})} + y^{(\text{train}) T} y^{(\text{train})}) = 0$$

$$w = (X^{(\text{train}) T} X^{(\text{train})})^{-1} X^{(\text{train}) T} y^{(\text{train})}$$

"normal equations"

extension to affine function:

$$\hat{y} = w^T x + b. \quad \text{bias parameter - without any input, output biased to } b. \quad \begin{array}{c} \| 1 \dots 1 \| \\ \vdots \end{array} \quad \begin{array}{c} 1 \dots 1 \\ \vdots \end{array} \quad \text{(not stratified)} \quad \begin{array}{c} 1 \dots 1 \\ \vdots \end{array}$$

$$y = w^T x + b$$

(note: can just add an extra entry to x that is always 1.) (not statistical bias)

so, typically call affine - "linear"

(5.2) Capacity, Overfitting, Underfitting

generalization - ability to perform well on previously unobserved inputs

training error: error on training data (can always find a model that makes this 0)

testing/generalization error: this is what separates ML from optimization

expected value of the error on new input

statistical learning theory

assumes data generating process (for both training and test data) follow

i.i.d. assumptions

\uparrow **identically distributed** - training & test sets drawn from same probability distribution

independent - each example in dataset is independent from others

shared underlying distribution, P_{data}

with i.i.d. assumption, randomly selected model: $\text{training error} = \frac{\text{expected}}{\text{test error}}$

but, with model learned on training data set: $\text{training error} \leq \text{expected test error}$

central challenges of ML:

1. make training error small \leftarrow avoid underfitting

2. make gap between training & test error small \leftarrow avoid overfitting.

capacity - ability to fit a wide variety of functions

high capacity \rightarrow overfitting memorize training set

low capacity \rightarrow underfitting

hypothesis space - set of functions a learning algorithm may select

e.g., linear regression: all linear functions

increase capacity \rightarrow all polynomials of degree 2

$$\hat{y} = b + w_1 x + w_2 x^2$$

note: still linear func. of parameters, so can use closed form lin. reg solution

ML algorithms perform best when capacity is appropriate for true complexity of the task

representational capacity / effective capacity

functions that can be represented

functions that can actually be learned

Occam's razor - principle of parsimony among competing hypotheses that match observations, pick simplest

Vapnik-Chervonenkis (VC) dimension - measure of capacity of binary classifier: largest m for which there exists a training set of m different x points that the classifier can label arbitrarily.

theoretical results in statistical learning theory

bound difference between trivial case & realistic

difference between training error & generalization error based on model (small then large)

size on model capacity), # of training examples
(upper) (lower)

bounds rarely useful for deep learning -

- difficult to determine capacity

effective capacity limited by optimization algorithm

little theoretical understanding of general nonconvex optimization algorithms.

• A family of general nonconvex optimizers used in deep learning

used in deep learning

nonparametric models — not limited to learning finite parameter vector

Nearest neighbor regression - store training set X, y , lookup closest

$$\hat{y} = y_i \text{ where } i = \arg \min \left\| X_{i,:} - x \right\|_2^2.$$

achieves minimum training error (if break ties)

wrapping parametric = outer loop that increases polynomial degree
for linear

Bayes error - error incurred by oracle making predictions from true distribution

(5.2.1) No Free Lunch Theorem (Wolpert 1996)

averaged over all possible data-generating distributions, every classification algorithm has the same error rate on previously unobserved points.

but - real data is not generated arbitrarily

(5.2.2) Regularization

NFL \Rightarrow must design ML algorithms to perform well on a specific task

idea: give preference to functions with certain properties

(will only prefer others, if they fit data much better)

e.g., weight decay for lin reg:

$$J(w) = MST + \lambda w^T w, \quad \text{minimize } J(w) \text{ to have}$$

$$J(w) = M S^T_{\text{train}} + \lambda w^T w$$

minimize $J(w)$ to have
preference for weights
parameter that controls with smaller L^2 norm
strength of preference

regularizer - penalty added to cost function

regularization - any modification made to a learning algorithm that is intended to reduce generalization error (but not training error).

philosophy of deep learning: a wide range of tasks (such as all the intellectual tasks that people can do) may all be solved effectively using very general-purpose forms of regularization.