The Last Few Bits.

CS771: Introduction to Machine Learning Pivush Rai

Debugging ML Algorithms



What is going wrong?

- What to do when our model (say logistic regression) isn't doing well on test data
 - Use more training examples?
 - Use a smaller number of features?
 - Introduce new features (can be combinations of existing features)?
 - Try tuning the regularization parameter?
 - Run (the iterative) optimizer longer, i.e., for more iterations
 - Change the optimization algorithm (e.g., GD to SGD or Newton..) or the learning rate?
 - Give up and switch to a different model (e.g., SVM or deep neural net)?



High-Bias or High-Variance?

- The bad performance (low accuracy on test data) of a model could be due either
 - High Bias: Too simple model; doesn't even do well on training data
 - High Variance: Even small changes in training data lead to high fluctuation in model's performance
- High bias means underfitting, high variance means overfitting
- Looking at the training and test error can tell which of the two is the case



Learning from Imbalanced Data



Learning when classes are imbalanced

When classes are imbalanced, even a "stupid" classifier can give high accuracy but looking at accuracy alone may be misleading



Solution 1: Balancing the training data

- Can balanced the training data by
 - Under-sampling the majority class examples



Weighted loss function with much larger importance given to loss function terms of positive examples than negative examples

$$\mathcal{L}(\boldsymbol{w}) = \sum_{i=1}^{N} \beta_{y_i} \ell(\boldsymbol{x}_i, y_i, \boldsymbol{w})$$
where $\beta_{+1} \gg \beta_{-1}$

cost/weights

33332

"negative" examples get weight 1 "positive" examples get a much larger weight change learning algorithm to optimize

Add costs/weights to the training set

weighted training error

Over-sampling the minority class examples





Equivalent to

99.997% not-phishing



Solution 2: Changing the loss function

Don't use loss functions that define loss or accuracy on per-example basis

This loss function is a simple sum of losses on individual training examples. Not ideal for imbalanced classes

$$L(\boldsymbol{w}) = \sum_{i=1}^{N} \ell(x_i, y_i, \boldsymbol{w})$$

- Instead, use loss function that use example pairs (one positive and one negative)
- Assuming our model to be defined by some function $f(\mathbf{x})$ (e.g., $\mathbf{w}^{\mathsf{T}}\mathbf{x}$), define a loss

An input with positive label
$$l(f(\mathbf{x}_n^+), f(\mathbf{x}_m^-)) = \begin{cases} 0, & \text{if } f(\mathbf{x}_n^+) > f(\mathbf{x}_m^-) \\ 1, & \text{otherwise} \end{cases}$$

Now we don't care about per-example accuracy but care about whether the positive examples get a higher score than the negative examples (i.e., we are only preserving their relative rank)

Such loss functions can known as "pairwise loss functions"

$$> \sum_{n=1}^{N_+} \sum_{m=1}^{N_-} \ell(f(\boldsymbol{x}_n^+), f(\boldsymbol{x}_m^-)) + \lambda R(f)$$



Ensemble Methods



Some Simple Ensembles

Voting or Averaging of predictions of multiple models trained on the same data



Stacking": Use predictions of multiple already trained models as "features" to train a new model and use the new model to make predictions on test data



Mixture of Experts (MoE) based Ensemble

- Mixture of Experts (MoE) is a very general idea
- We assume *m* "simple" models, usually of the same type, e.g., *m* linear SVMs or *m* logistic regression models, or *m* deep neural nets (usually all with same architecture)



Ensembles using **Bagging** and **Boosting**

- Both use a single training set ${\cal D}$ to learn an ensemble consisting of several models
- Both construct N datasets from the original training set $\mathcal D$ and learn N models



Bagging

Boosting

CS771: Intro to ML

- Bagging can do this in <u>parallel</u> for all the N models
- Boosting requires a <u>sequential</u> approach for N rounds

Figure credit: https://www.blog.dailydoseofds.com/p/an-animated-guide-to-bagging-and

Bagging (Bootstrap Aggregation)



Figure adapted from: https://www.blog.dailydoseofds.com/p/an-animated-guide-to-bagging-and



Figure adapted from: https://www.blog.dailydoseofds.com/p/an-animated-guide-to-bagging-and

A Boosting Algo: AdaBoost (Adaptive Boosting)

Importance of the training example (x_i, y_i)

or estimate it during training

importance for all training examples



- AdaBoost is based on optimizing such a loss function
 - Initialize the ensemble as $\mathcal{E} = \{\}$ and $\boldsymbol{\beta}$ as $\boldsymbol{\beta}^{(0)} = \left[\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}\right]$

• For round
$$t = 1, 2, ..., T$$

Fin

- $w^{(t)} = \operatorname{argmin}_{w} \sum_{i=1}^{N} \beta_{i}^{(t-1)} \ell(\mathbf{x}_{i}, \mathbf{y}_{i}, \mathbf{w})$ and add it to ensemble $\mathcal{E} = \{\mathcal{E} \cup \mathbf{w}^{(t)}\}$
- Define the total loss of $w^{(t)}$ as $L(w^{(t)}) = \sum_{i=1}^{N} \beta_i^{(t-1)} \ell(x_i, y_i, w^{(t)})$ Or the importance weighted total error
- Compute the "importance" of $w^{(t)}$ for the \mathcal{E} as $\alpha_t = f(L(w^t)) \int_{\text{total loss } L(w^t)} f$ is some function such that α_t is high if total loss $L(w^t)$ is low, and vice-versa
- Increase/decrease importance β_i of each training instance (x_i, y_i) for next round as

$$\beta_{i}^{(t)} \propto \begin{cases} \beta_{i}^{(t-1)} \times \exp\left(\alpha_{t}\ell(\boldsymbol{x}_{i}, y_{i}, \boldsymbol{w}^{(t)})\right) & (\text{Increase if } \boldsymbol{w}^{(t)} \text{ mispredicted } (\boldsymbol{x}_{i}, y_{i})) \\ \beta_{i}^{(t-1)} \times \exp\left(-\alpha_{t}\ell(\boldsymbol{x}_{i}, y_{i}, \boldsymbol{w}^{(t)})\right) & (\text{Decrease if } \boldsymbol{w}^{(t)} \text{ correctly predicted on } (\boldsymbol{x}_{i}, y_{i})) \end{cases}$$
al model is $\hat{\boldsymbol{w}} = \sum_{t=1}^{T} \alpha_{t} \boldsymbol{w}^{(t)}$ importance-weighted average of all $\boldsymbol{w}^{(t)}$'s CS771: Intro to ML

AdaBoost: An Illustration

- Suppose we have a binary classification problems with each input having 2 features.
- Suppose we have a weak model like a simple DT (decision stump)
- Illustration of AdaBoost using a decision stump if run for 3 rounds
 Round 1
 Round 2



- The ensemble represents the overall model
- We got a nonlinear model from 3 simple linear models
- Note that the ensemble was constructed sequentially



Original dataset

Gradient Boosting

- Consider learning a function f(x) by minimizing a squared loss $\frac{1}{2}(y f(x))^2$
- Gradient boosting is a sequential way to construct such f(x)
- For simplicity, assume we start with $f_0(x) = \frac{1}{N} \sum_{i=1}^N y_i$
- Given previously learned model $f_m(x)$, let's assume the following "improvement" to it

$$f_{m+1}(x) = f_m(x) + h(x)$$

Based on sequentially

- Thus the goal for the next round is to learn the "residual" $h(x) = y f_m(x)$
- Residual is negative gradient of the loss w.r.t. f(x) thus called "gradient boosting"
- The final model $f_M(x)$, once the residual is sufficiently small, is what we will use
- The idea of gradient boosting is applicable to classification too
- XGBoost (eXtreme Gradient Boosting) is a very popular grad boosting algo

Active Learning

(an example of learning with human-in-the-loop)



18

Active Learning

- Standard supervised learning is "passive" (learner has no control; we just give it data)
- We take a random sample of inputs, get them labelled by an expert, and train a model





Active Learning

In active learning, learner can request what training examples it wants to learn from



Learning in the wild



21

Domain Adaptation

- We may have a "source" model trained on data from some domain
- We might want to deploy it in a new domain
- Performance of the source model will suffer
- To prevent this, we usually perform "domain adaptation" or "transfer learning"
- These are broad terms covering a variety of techniques that "finetune" the source model using labelled/unlabeled data from the new domain



We do expect some

"commonality" (e.g., some common set of features)



The ending note..

- Good features are important for learning well
- The "classical" ML methods we studied in this course still continue to have high relevance
- Success of deep learning is largely attributed to (automatically learned) good features
- Deep learning is not a panacea often simple classical models can do comparably/better
- First understand your data (plot/visualize/look at some statistics of the data, etc)
- Always start with a simple model that you understand well
 - Try to first understand if your data really needs a complex model
- Think carefully about your features, how you compute similarities, etc.
- Helps to learn to first diagnose a learning algorithm rather than trying new ones
 - Understanding of optimization algos, loss function, bias-variance trade-offs, etc is important
- No free lunch. No learning algorithm is "universally" good

