

Classification in Big Data Environments

Clemens Damke

Intelligent Systems and Machine Learning Group (ISG)
Heinz Nixdorf Institute
Paderborn University
Warburger Straße 100
33098 Paderborn

Abstract. *TODO*

Keywords: Big Data · Automated Machine Learning

1 Introduction

Over the last few years Big Data processing has become increasingly important in many domains. This increase in the data volume also poses new challenges for machine learning applications. The training time of learners is usually polynomially dependent on the size of the training dataset \mathcal{D}_{train} , i. e. $\mathcal{O}(|\mathcal{D}_{train}|^\alpha)$, $\alpha \geq 1$. Since training has to be repeated for every iteration of cross-validation and hyperparameter search, always using the entire dataset quickly becomes infeasible. This paper gives an overview of approaches to tackle this problem with a focus on the domain of classification problems.

The process of finding a hypothesis can in general be split into three phases:

1. **Model selection:** At first a model has to be selected that determines the class of hypothesis spaces out of which the final hypothesis will be selected. The choice of model is often implicitly encoded in a learner L . Automating this step is non-trivial. In practice the model is typically selected by experts with domain specific knowledge about the problem at hand.
2. **Hyperparameter search:** Optimizing a vector λ in the hyperparameter space Λ_L of the learner L representing a hypothesis space \mathcal{H}_λ . A naïve approach to do this is to systematically try configurations using a grid search or a random search over Λ_L . To evaluate the quality of a given λ , L is usually trained on a training dataset \mathcal{D}_{train} using λ . This yields a hypothesis $\hat{h}_\lambda \in \mathcal{H}_\lambda$ that is evaluated using a validation dataset \mathcal{D}_{valid} . The goal of hyperparameter optimization is to minimize the loss $l(\lambda)$ of \hat{h}_λ on \mathcal{D}_{valid} , i. e. to find an approximation $\hat{\lambda}$ of $\lambda^* := \arg \min_\lambda l(\lambda)$.
3. **Training or parameter search:** Let w be a vector in the parameter space $W_{\mathcal{H}_\lambda}$, describing a hypothesis $h_{\lambda,w} \in \mathcal{H}_\lambda$ given a hyperparameter configuration λ . The goal of parameter search is to find an approximation \hat{h}_λ of the hypothesis $h_\lambda^* := \arg \min_{h_{\lambda,w}} \ell(\mathcal{D}_{train}|h_{\lambda,w})$, with $\ell(\mathcal{D}_{train}|h_{\lambda,w})$ being the empirical loss of $h_{\lambda,w}$ on a given training dataset \mathcal{D}_{train} according to some loss function ℓ . Depending on the learner L , various kinds of optimization methods are used to find this minimum, e. g. Bayesian optimization, quadratic programming or, if $\nabla_w \ell(\mathcal{D}_{train}|h_{\lambda,w})$ is

computable, gradient descent. The quality l of \hat{h}_λ is measured by the loss on a validation or test dataset, i. e. $l(\lambda) := \ell(\mathcal{D}_{\text{valid}}|\hat{h}_\lambda)$.

This paper is structured according to the last two phases, i. e. we will assume that the learner L is given. Section 2 describes ways to speed up the hyperparameter search. Section 3 then describes how to improve the training methods of existing learners. Most of the techniques described in this paper improve upon independent components of the hypothesis finding process allowing them to be combined.

2 Hyperparameter optimization

As described in the introduction, the goal of hyperparameter optimization is to find a global minimum of l . Since l is generally unknown, analytical methods or gradient descent cannot usually be applied. The only way to get information about l is to evaluate it, which is costly. There are multiple ways to reduce the total cost of those evaluations:

1. **Number T of evaluations of l :** During optimization multiple hyperparameter configurations $\lambda_1, \dots, \lambda_T$ will be evaluated using l . T is usually fixed when using a grid search or a random search. After evaluating T configurations, the best one is chosen. Those naïve approaches assume that $l(\lambda)$ is independent of $l(\lambda')$ for all pairs $\lambda \neq \lambda'$. We will see that this strong assumption of independence is not necessarily true which in turn allows reducing T .
2. **Training dataset size S :** The performance of a given configuration $l(\lambda)$ is computed by training the learner on $\mathcal{D}_{\text{train}}$ which is expensive for big datasets. By training on S instead of $|\mathcal{D}_{\text{train}}|$ datapoints the evaluation can be sped up.
3. **Number of training iterations E :** Depending on the learner, training often is an iterative process, e. g. gradient descent. To speed up hyperparameter optimization training could be terminated before convergence.

2.1 FABOLAS

The first approach we will discuss is called Fabolas (Fast Bayesian Optimization of Machine Learning Hyperparameters on Large Datasets) [4]. It can be applied to any learner L and is based upon two main ideas:

1. The validation loss l is modeled as a *Gaussian process* (GP) f based on the assumption that two configurations λ and λ' will perform similar if they are similar according to some kernel $k(\lambda, \lambda')$. The Gaussian process f is used as a surrogate to estimate the expected value and variance of l given λ . Using *Bayesian optimization* l will be probed at promising positions to iteratively improve f . Hyperparameter configurations that are expected to perform worse than the current optimum will not be probed. This effectively reduces T .
2. The training dataset size S is modeled as an additional hyperparameter of f giving the optimizer an additional degree of freedom. This allows extrapolating the value of l when trained on the complete dataset while only probing smaller subsets which effectively reduces S .

We will now describe how those two ideas can be applied.

Gaussian processes A Gaussian process is a family of *random variables* (RVs) $(X_\theta)_{\theta \in \Theta}$, s. t. every finite subset of them follows a multivariate normal distribution. More intuitively it can be understood as a probability distribution over functions $f : \Theta \rightarrow \mathbb{R}$ where $X_\theta \triangleq f(\theta)$. Prior knowledge about the likelihood of each f is described by a prior mean function $\mu_0(\theta) = \mathbb{E}[f(\theta)]$ and a positive-definite kernel $k(\theta, \theta') = \text{Cov}(f(\theta), f(\theta'))$. The covariance kernel models how informative it is to know $f(\theta)$ to determine $f(\theta')$.

Let $\mathcal{D}_n = \{(\theta_i, \mathbf{y}_i)\}_{i=1}^n$ denote a set of observations. Those observations can be used to update the means and variances of the RVs via GP regression. This collapses the space of possible functions f to those functions that align with \mathcal{D}_n (see fig. 2):

$$\begin{aligned} \mathbf{m} &:= (\mu_0(\theta_1), \dots, \mu_0(\theta_n))^T \\ \mathbf{k}(\theta) &:= (k(\theta_1, \theta), \dots, k(\theta_n, \theta))^T \\ \mathbf{K} &\in \mathbb{R}^{n \times n}, \mathbf{K}_{ij} := k(\theta_i, \theta_j) \\ \mathbb{E}[f(\theta) | \mathcal{D}_n] &:= \mu_n(\theta) = m_0(\theta) + \mathbf{k}(\theta)^T \mathbf{K}^{-1}(\mathbf{y} - \mathbf{m}) \end{aligned} \quad (1)$$

$$\text{Cov}(f(\theta), f(\theta') | \mathcal{D}_n) := k(\theta, \theta') - \mathbf{k}(\theta)^T \mathbf{K}^{-1} \mathbf{k}(\theta') \quad (2)$$

Fabolas works by modeling the loss function l as a Gaussian process $f \sim \mathcal{GP}(m, k)$ with parameter set $\Theta := \Lambda \times [0, 1]$ where $\mu_0(\lambda, s) = \mathbb{E}[f(\lambda, s)] = \mathbb{E}[l(\lambda) | \text{training size } s]$. To model the covariances between different combinations of hyperparameters and training set sizes, the following product kernel is used:

$$k((\lambda, s), (\lambda', s')) := k_{\text{MATÉRN5}}(d_M(\lambda, \lambda')) \cdot k_{\text{lin}}(s, s') \quad (3)$$

Here $k_{\text{MATÉRN5}}$ denotes the stationary Matérn kernel ($\nu = 5/2$) with d_M being the Mahalanobis distance between the two compared hyperparameter configurations. k_{lin} essentially is a simple linear kernel modeling the assumption that l monotonically decreases when s is increased. We will only give an intuition for this choice of kernel and refer to Klein et al. [4] for the details.

The Mahalanobis distance d_M is used instead of the Euclidean distance because the hyperparameters in a configuration typically use very different scales and are in some cases also correlated. Figure 1 gives an intuition for this.

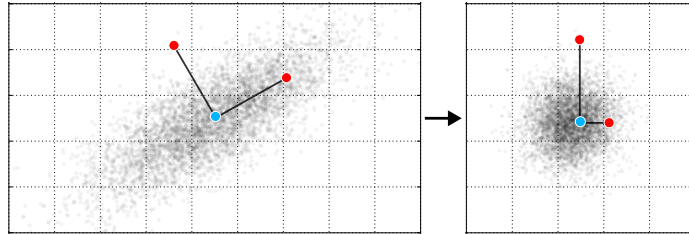


Fig. 1: Intuition for the Mahalanobis distance. Using the Euclidean distance the red points would be equally far away from the blue one. The Mahalanobis distance fixes this by first normalizing the hyperparameters and removing correlations.

Based on the Mahalanobis distance between two configurations λ, λ' the MATÉRN5 kernel is used to compute a covariance. The class of Matérn kernels interpolates between the Gaussian (SQ-EXP) and the exponential (MATÉRN1) kernel (see fig. 2). Because the exponential kernel drops off quickly, configurations quickly become uncorrelated which causes noisy samples. The Gaussian kernel drops off less quickly causing smoother samples. Fabolas uses MATÉRN5 as it empirically fits the smoothness of typical loss functions l quite well.

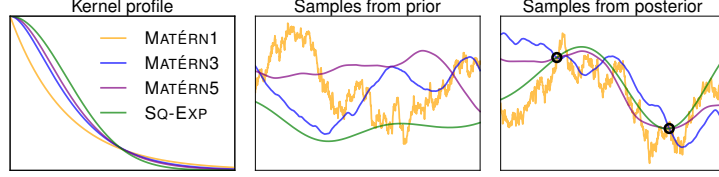


Fig. 2: Comparison between different covariance kernels. The middle shows randomly sampled functions f using different kernels. The right shows random samples after two f values were observed and incorporated into the model via GP regression.

Bayesian optimization To find $\arg \min_{\lambda} l(\lambda)$ the bias and variance of f has to be reduced by probing l at promising positions. This is called Bayesian optimization. The estimated minimum after n probes is described by $\arg \min_{\lambda} \mu_n(\lambda, s = 1)$, i. e. the configuration with the smallest predicted error on the full test dataset. To reduce the number of probes required until this minimum converges, an *acquisition function* is used. Its role is to trade-off exploration vs. exploitation of l by describing the expected utility of probing (λ_{n+1}, s_{n+1}) given a set of previous probes \mathcal{D}_n . Fabolas uses an acquisition function that rates configurations by their *information gain* per computation time:

$$a_F(\lambda, s) := \frac{1}{c(\lambda, s)} \mathbb{E}_y \left[p(y | \lambda, s, \mathcal{D}_n) \cdot \text{KL}_{\hat{\lambda}}(p_{\min}(\hat{\lambda} | \mathcal{D}_n \cup \{(\lambda, s, y)\}) || u(\hat{\lambda})) \right] \quad (4)$$

$$p_{\min}(\lambda | \mathcal{D}) := p(\lambda \in \arg \min_{\lambda'} f(\lambda', s = 1) | \mathcal{D})$$

It measures the expected amount of available information about the optimal configuration if a given configuration were probed, i. e. the Kullback-Leibler divergence between the density p_{\min} of λ being optimal after a probe and the uniform density u . This information gain of a probe is compared to its expected associated computation time c . c is estimated using a separate Gaussian process that is maintained alongside f . Fabolas considers the cost of a probe because it tries to minimize the total optimization time not the total number of probes.

Since it is infeasible to compute a_F numerically its maximum is estimated using *Markov-Chain Monte Carlo* (MCMC). The estimated most promising configuration will be probed. The resulting loss value and runtime are then used to update the loss model f and cost model c via GP regression.

Evaluation Fabolas was evaluated in *support vector machine* (SVM) and *convolutional neural network* (CNN) optimization tasks on the MNIST and CIFAR-10 dataset respectively. Figure 3 compares Fabolas to the following other hyperparameter optimization approaches:

- **Random Search:** Simple random hyperparameter search. Each configuration is evaluated on the full dataset.
- **Entropy Search & Expected Improvement:** Bayesian optimization methods that always evaluate on the full dataset. Expected Improvement uses an acquisition function that simply probes at the current expected optimum. Entropy Search uses an acquisition function similar to the one used by Fabolas but without the cost model.
- **MTBO- N (Multi-Task Bayesian Optimization):** Like Fabolas but restricts probes to two sizes $s \in \{1/N, 1\}$, i. e. either a small subsample or the entire dataset is used. Multiple values for N were evaluated: 4, 32 and 512.

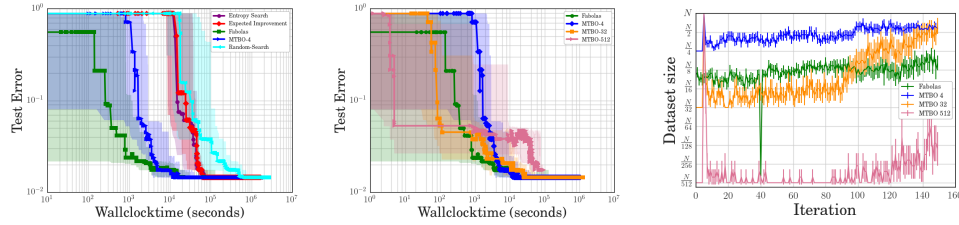


Fig. 3: SVM optimization on the MNIST dataset. (Left) Comparison of the test performance over time of different optimizers. (Middle) Comparison of Fabolas with different MTBO subsample sizes. (Right) Comparison of the subsample sizes s that MTBO and Fabolas choose for their probes. The average over 10 runs is depicted.

All Bayesian optimization approaches are at least one order of magnitude faster than random search. By allowing two probing sizes, MTBO is one additional order of magnitude faster. Depending on the choice of N MTBO sometimes improves faster than Fabolas initially. Once Fabolas starts improving it does however find a good configuration about one order of magnitude faster than MTBO. The optimal configuration is found at roughly the same time by both Fabolas and MTBO. Overall Fabolas finds a good configuration between 100 and 1000 times faster than random search does. Similar results are obtained when optimizing CNNs on CIFAR-10.

2.2 Learning Curve Extrapolation

The second approach for speeding up hyperparameter optimization focuses on reducing the number of training iterations E . It can in principle be applied to any iterative learner and can be integrated into any hyperparameter optimizer. The idea is to monitor the learning curve of a learner during training with a hyperparameter configuration λ . If it is unlikely that a good accuracy will be reached with λ , training will be terminated before convergence.

The method was first described by Domhan et al. [3] in the context of hyperparameter optimization for *deep neural networks* (DNNs) that are trained using *stochastic gradient descent* (SGD).

Extrapolation method Let $y_{1:n}$ denote the observed learning curve of SGD after n iterations, i.e. the sequence of training accuracies $y_i \in [0, 1]$. Normally SGD iterations would be run for each hyperparameter configuration λ until convergence or until a maximum number of iterations E has been reached. The learning curve extrapolation optimization works by predicting y_E every p iterations:

```

1:  $\hat{y} \leftarrow -\infty$ 
2: for  $\lambda \leftarrow$  next hyperparameter configuration to evaluate do
3:    $n \leftarrow 0$ 
4:   repeat
5:     Run  $p$  SGD iterations using  $\lambda$  with resulting accuracies  $y_{(n+1):(n+p)}$ .
6:      $n \leftarrow n + p$ 
7:     Estimate  $P(y_E < \hat{y} \mid y_{1:n})$ .
8:   until SGD converged  $\vee n \geq E \vee P(y_E < \hat{y} \mid y_{1:n}) > \delta$ 
9:   if  $y_n > \hat{y}$  then  $\hat{y} \leftarrow y_n$  end if
10: end for

```

The prediction step (line 7) uses a probabilistic model. Similar to Fabolas, a distribution over candidate functions is fit to the observations $y_{1:n}$. Unlike Fabolas however, which uses a flexible non-parametric GP model, we use prior knowledge about the shape of learning curves to restrict the model to parameterized, increasing, saturating functions. More specifically, the learning curve $y_{1:n}$ is modeled as a linear combination f_{comb} of a family of given functions.

$$f_{comb}(t \mid \xi) := \sum_{k=1}^K w_k f_k(t \mid \theta_k), \quad \xi = (w_1, \dots, w_K, \theta_1, \dots, \theta_K, \sigma^2) \quad (5)$$

$$y_t \sim \mathcal{N}(f_{comb}(t \mid \xi), \sigma^2) \quad (6)$$

Domhan et al. [3] use $K = 11$ types of functions $\{f_1, \dots, f_K\}$ that are each parameterized by $\{\theta_1, \dots, \theta_K\}$. The assumption is that every function type captures certain aspects of learning curves. By allowing linear combinations a more powerful model can be obtained. Figure 4 illustrates this idea. To estimate the probability $P(y_E < \hat{y} \mid y_{1:n})$ MCMC is used to sample S learning curves $\{\xi_1, \dots, \xi_S\}$ from the posterior

$$P(\xi \mid y_{1:n}) \propto P(y_{1:n} \mid \xi) P(\xi) \quad (7)$$

$$P(y_{1:n} \mid \xi) = \prod_{t=1}^n \mathcal{N}(y_t; f_{comb}(t \mid \xi), \sigma^2) \quad (8)$$

$$P(\xi) \propto \mathbb{1}[f_{comb}(1 \mid \xi) < f_{comb}(E \mid \xi) \wedge \forall k : w_k > 0] \quad (9)$$

The prior on ξ is used to model the fact that learning curves do not typically decrease over time. Given the learning curve samples, we can now estimate

$$\begin{aligned}
P(y_E < \hat{y} \mid y_{1:n}) &= \int P(\xi \mid y_{1:n}) P(y_E < \hat{y} \mid \xi) d\xi \\
&\approx \frac{1}{S} \sum_{s=1}^S \Phi(\hat{y}; f_{comb}(E \mid \xi_s), \sigma^2)
\end{aligned} \quad (10)$$

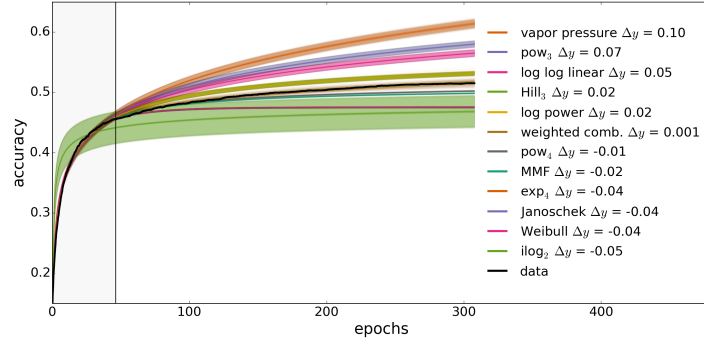


Fig. 4: Comparison of an observed learning curve (black) with the 11 types of learning curve models and a linear combination of them. Each type is parameterized to fit the first 50 observations $y_{1:50}$. As can be seen in the legend on the left, the linear combination has the smallest deviation Δy from the observed data after 300 iterations.

Evaluation The early termination method we just described was evaluated on the CIFAR-10, CIFAR-100 and MNIST dataset. Figure 5 shows the behavior of early termination and the obtained speedup on CIFAR-10. As expected, configurations with learning curves that tend to approach low accuracies are terminated early. Configurations with high accuracies are run until convergence. This approach is consistently speeds up the hyperparameter optimization by a factor of two across the tested datasets while reaching the same quality.

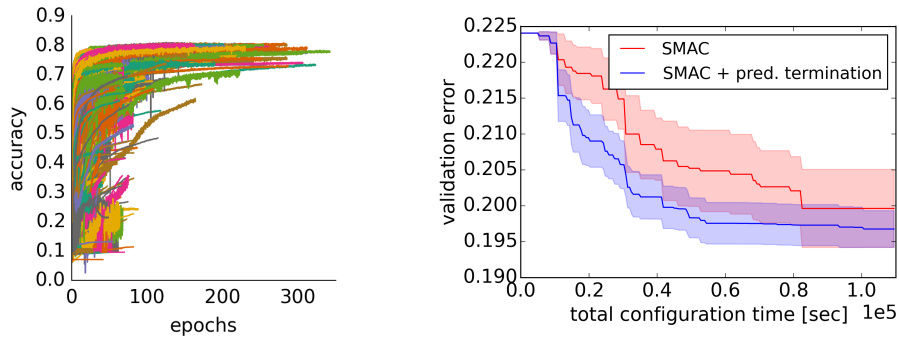


Fig. 5: Evaluation on the CIFAR-10 dataset. The left graph shows the learning curves of all hyperparameter configurations that were evaluated. The right graph shows the average validation error over time.

3 Optimizing Training

Now follows an overview of approaches to speed up the training process. We will discuss four approaches:

1. A general purpose method that combines subsampling with bootstrapping.

2. An iterative method to select the optimal subsample size during gradient descent.
3. Improving the quality of subsampling for logistic regression by weighing the samples.
4. Speeding up the training of SVMs via k -means clustering.

3.1 Bag of Little Bootstraps

The first approach we will discuss is called *Bag of Little Bootstraps* (BLB) [5]. It combines subsampling with bootstrapping and is particularly well suited for parallelized implementations.

In the context of Big Data training typically cannot be performed on the entire dataset. A naïve way to solve this problem is to simply train on a random b out of n subsample of the data $\mathcal{D}_{train} = \{X_1, \dots, X_n\}$. This approach is highly sensitive to noise in the training dataset, especially if $b \ll n$. To overcome this problem bootstrapping can be used. The regular n out of n bootstrapping technique for variance reduction is not suitable for big datasets because it uses 63% of the training data on average. However the b out of n bootstrapping approach can in principle be applied. It uses s samples $\{\tilde{X}^{(i)} = (\tilde{X}_1^{(i)}, \dots, \tilde{X}_b^{(i)}) \mid 1 \leq i \leq s\}$. Since this approach independently learns s hypotheses \tilde{h}_i on small datasets $\tilde{X}^{(i)}$, their parameterizations $\tilde{\theta}_i$ tend to have large confidence intervals. Because of that, the quality of the combined hypothesis is strongly dependent on b . BLB reduces this dependence.

Intuition BLB is a simple extension of the b out of n bootstrap that is consistently more robust regarding the choice b across datasets. The basic idea is to add another sampling step. BLB uses each subsample $\tilde{X}^{(i)}$ as a seed for n out of b sampling. This yields samples $\{X^{*(i,k)} = (X_1^{*(i,k)}, \dots, X_n^{*(i,k)}) \mid 1 \leq i \leq s, 1 \leq k \leq r\}$ that each contain at most b different elements. Training is then run on the big samples X^* instead of the small seed samples \tilde{X} .

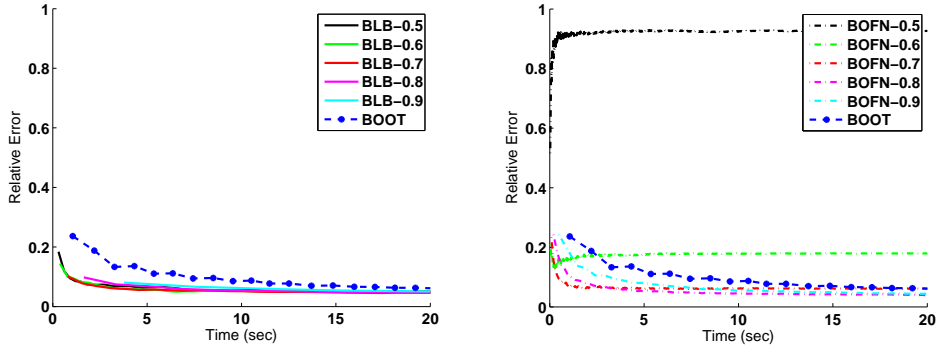


Fig. 6

Evaluation

3.2 Subsample Size Selection for Gradient Descent

Next we will discuss gradient descent based optimization. The sample size S that describes the number of datapoints that are considered in a single gradient descent step heavily influences the optimizer's behavior:

- In the stochastic approximation regime small samples, typically $S = 1$, are used. This causes fast but noisy steps.
- In the batch regime large samples, typically $S = |\mathcal{D}_{train}|$, are used. Steps are expensive to compute but more reliable.

Both extremes are usually not suitable for Big Data applications. Very small samples cannot be parallelized well, making them a bad fit for the compute clusters that are typically available nowadays. The gradients for very big samples however are often too slow to compute. S should ideally lie somewhere in between.

Byrd et al. [2] describe an iterative algorithm that dynamically increases S as long as this promises to significantly reduce the gradient noise.

3.3 Subsampling for Logistic Regression

Subsampling usually increases the mean squared error (MSE) of the resulting hypothesis compared to one that is trained on the full dataset. OSMAC [6] is a method that improves upon naïve subsampling by weighing the samples.

TODO: explain OSMAC

3.4 SVM-KM

To speed up the training of SVMs Almeida et al. [1] proposed a simple method that reduces the dataset size via k -means clustering. It can be described as a three-step procedure:

1. Group the training samples \mathcal{D}_{train} into k clusters C_1, \dots, C_k with centers c_1, \dots, c_k where k should be determined via hyperparameter optimization.
2. Check for each cluster C_i whether all associated datapoints belong to the same class, i. e. $\exists z \in \{+1, -1\} : \forall (x, y) \in C_i : y = z$. If yes, all datapoints in C_i are removed from \mathcal{D}_{train} and replaced by c_i . If not, they are kept in the dataset. The intuition behind this is that clusters with points from multiple classes might be near the decision boundary so they are kept to serve as potential support vectors.
3. Finally standard SVM training is performed on the reduced training dataset.

Evaluation *TODO*

4 Conclusion

TODO

References

- [1]M. Barros de Almeida, A. de Padua Braga, and J. P. Braga. “SVM-KM: speeding SVMs learning with a priori cluster selection and k-means”. In: *Proc. Vol.1. Sixth Brazilian Symp. Neural Networks*. Nov. 2000, pp. 162–167 (cit. on p. 9).
- [2]Richard H. Byrd, Gillian M. Chin, Jorge Nocedal, and Yuchen Wu. “Sample size selection in optimization methods for machine learning”. In: *Mathematical Programming* 134.1 (2012), pp. 127–155 (cit. on p. 9).
- [3]Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. “Speeding Up Automatic Hyperparameter Optimization of Deep Neural Networks by Extrapolation of Learning Curves”. In: *Proceedings of the 24th International Conference on Artificial Intelligence. IJCAI’15*. Buenos Aires, Argentina: AAAI Press, 2015, pp. 3460–3468 (cit. on p. 6).
- [4]Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. “Fast Bayesian Optimization of Machine Learning Hyperparameters on Large Datasets”. In: (May 23, 2016). arXiv: 1605.07079v2 [cs.LG] (cit. on pp. 2, 3).
- [5]Ariel Kleiner, Ameet Talwalkar, Purnamrita Sarkar, and Michael I. Jordan. “A Scalable Bootstrap for Massive Data”. In: (Dec. 21, 2011). arXiv: 1112.5016v2 [stat.ME] (cit. on p. 8).
- [6]HaiYing Wang, Rong Zhu, and Ping Ma. “Optimal Subsampling for Large Sample Logistic Regression”. In: (Feb. 3, 2017). arXiv: 1702.01166v2 [stat.CO] (cit. on p. 9).