# Abstract

I'm Abstract, I'm Abstract, I'm Abstract,I'm Abstract,I'm Abstract,I'm Abstract

**Keywords**: words, words, words

# Introduction

**Large data sizes and dangers+cost of collecting data in one place** Mention figure 1 and how ideally that's how we would like to do things

**Advances in machine learning and advantages of using more data**



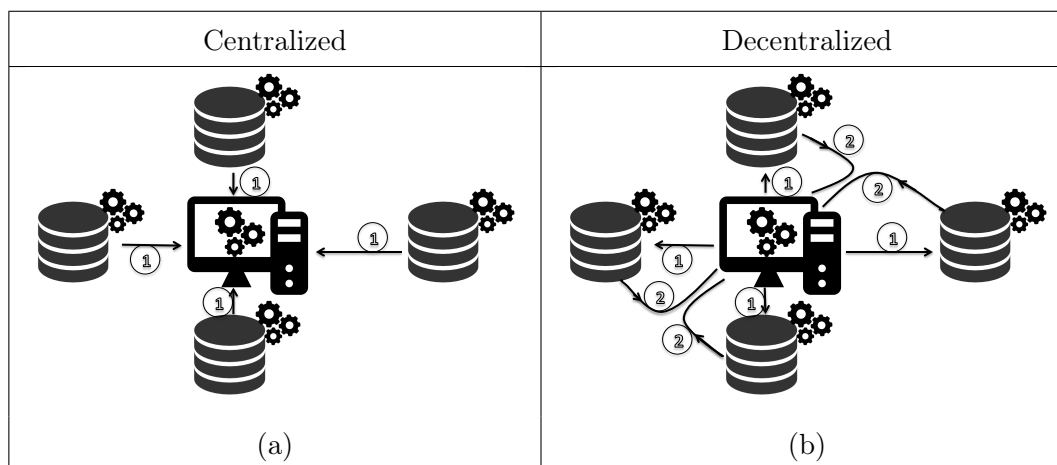| Centralized | Decentralized |
| --- | --- |
| (a) | (b) |

Table 1: (a)The simplest method of data-sharing would involve all DOs sending their data to a single location. That location would then execute the relevant machine learning algorithms. (b) In our model of decentralized scheme, one central hub can communicate with all the DO's without directly requiring the data to be sent to the central hub. The central hub coordinates the execution of the relevant machine learning algorithm and through a potentially iterative process, a single set of parameters is estimated for the machine learning algorithm.

**Introduction to GLM and their extensive use**

# Methods

In a broad sense, all regression based classifiers are consisted of three parts: representation, evaluation and optimization (4). The **representation** generally describes the relationship between parameters and limits the class of models that the learner can represent. The **evaluation function** (also known as cost function) is a tool for evaluating different parameter settings of the model. Finally, the **optimization technique** is the algorithm used to find the optimum model-parameters with respect to the chosen evaluation function. In this work, we will focus on generalized linear

| Type of GLM | Mean function $E[\mathbb{Y}]$ | Semantics |
|---|---|---|
| Linear regression | $\mu = X\beta$ | Outcome is a linear function of the parameters with range $(-\infty, \infty)$ |
| Logistic regression | $\mu = \frac{\exp(X\beta)}{1+\exp(X\beta)}$ | Probability of a categorical outcome based on a binomial distribution |
| Poisson regression | $\mu = \exp X\beta$ | Modeling Poisson distributed $\frac{\text{counts}}{\text{time interval}}$ |
| Multinomial regression | $\mu = \frac{\exp(X\beta)}{1+\exp(X\beta)}$ | Probability of each outcome for a K-sided loaded die |

Table 2: List of most commonly used link functions, the GLM regression they correspond to and their common use case.

16 models (GLMs). We will first briefly review the model and the cost function and we will focus the
17 bulk of this section on reviewing the optimization methods and how these methods can be optimized
18 in a decentralized way.

**GLM representation**: GLMs treat each response variable as a random draw from a distribution in the exponential family where the mean of the distribution can be related to a linear function of the covariates (17). More concretely, A GLM models the expected value of the outcome variable ($\mathbb{Y}$) as:

$$E(\mathbb{Y}) = \mu = g^{-1}(X\beta) \tag{1}$$

19 Where $X$ is a matrix with covariates for each individual as its rows. $\beta$ is the vector of unknown
20 parameters and g is a link function which is assumed to be known. In this setting, the link function
21 ($g$) provides flexibility to model various types of outcomes. Table 2, provides a short list of examples
22 of common link functions and their typical use cases [perhaps cite a review here].

23 **GLM evaluation function**: The unknown parameters ($\beta$) are chosen to minimize a particular
24 cost function. There are a few natural ways to define a cost function for GLMs. Amongst these
25 various methods [cite bayesian approach and least square fits to variance stabilized response related
26 papers]), using maximum likelihood to define a cost function remains one of the most popular
27 approach. Here we will only consider optimization techniques for such cost functions.

28 **GLM optimization techniques** vary depending on the details of the setting and will constitute
29 the main focus of this review. While, the default for some common packages (glm2 in R for example
30 (15)) uses iteratively reweighted least squares (reviewed by [cite]), as explained in the introduction,
31 these methods may not be practical if the data is very high dimensional or, otherwise, cannot be

gathered in one location. In these settings, many packages (such as glmnet in R, (7), Vowpal Wabbit (11) and Python's scikit-learn) use sub-gradient based methods (reviewed in section[TODO]) to fit the model parameters.

############### ABOVE HAS BEEN UPDATED + TWO NEW FIGURES BE-LOW

Unsurprisingly, much work has been done for

Furthermore, in this setting the updates are expensive when new data becomes available. In recent years, much work has been done on fitting GLM's under the online setting [cite cite cite ] or when the data is distributed. In this review, we focus our discussion along the axis of table [number]. In particular, we will discuss methods that can be applied when all the data is available (offline) on one machine or distributed amongst numerous data owners (DO) as well as when the data is coming in and one hopes to efficiently update the previous estimates with the information from the newly available data.

## Centralized vs. distributed

An important factor in choosing the optimization technique is whether the data can be centralized and loaded into the memory. Here, we will consider a model where a central hub coordinates many DOs and solves the problem without copying the data from the DO to the hub (fig**??**) We will assume that each DO has compute power but cannot hold the entire dataset in memory. The general, the algorithm starts with 1) the DO informing the hub of arrival of new sample(s), 2) the hub proceeds to send the DO relevant information and 3) the DO sends information back on how the parameters should be updated.

### Stochastic Gradient Descent

Before describing the algorithms for each of the scenarios in table 1, we briefly review the stochastic gradient descent (SGD) algorithm. For independent samples, the MLE-derived cost function for GLM family ($C(\beta; X, y)$) can separated into the sum of cost functions over all the individuals.

$$C(\beta; X, y) = \sum_{i=1}^{N} C_i(\beta) \tag{2}$$

Where $N$ is the number of individuals and $C_i$ is the cost function calculated over the $i$th data

point. Then the gradient descent update rule is given as

$$\beta_{(t+1)} = \beta_t - \eta_t \sum_{i=1}^{N} \nabla C_i(\beta_t) \qquad [3]$$

54 Where $\nabla C_i(\beta^t)$ is the gradient of the cost function evaluated at the previous estimate of $\beta$, $\eta_t$ is

55 the step size for the $t$th iteration of the algorithm. For large $N$ computing the sum can become

56 computationally expensive particularly given the relative slow rate of convergence for this algorithm

57 (linear convergence for a $\beta$-smooth and $\alpha$-convex function **check these conditions and perhaps**

58 **replace that with GLM family for simplicity (if it is true)**). (Batch) stochastic gradient

59 descent (SGD) attempts to circumvent this problem by updating the estimate based on the gradient

60 from a randomly chosen data points (or batch). SGD trades off faster convergence rate of GD with

61 the smaller cost of each individual updates and for the same set of functions (**match this to GD's**

62 **claim**), converges in time independent of the total number of samples present (cite Optimization

63 Methods for Large-Scale Machine Learning).

The choice of $\eta$ is crucial for convergence of SGD. In particular, convergence requires:

$$\sum_{t=1}^{\infty} \eta_t^2 < \infty; \quad \sum_{t=1}^{\infty} \eta_t = \infty \qquad [4]$$

64 Which intuitively means that the step-size should be chosen to be small but not too small. A

65 learning schedule with $\eta_t = \frac{c_1}{(c_2+t)^\gamma}$ and $0.5 < \gamma \leq 1$ satisfies this condition. However, it must be

66 noted that still the values of $c_1, c_2 > 0$ need to be chosen judiciously. In particular, if the step-sizes

67 are too large, the algorithm may not converge and if the step-sizes are too small, the convergence

68 rate may be too small (cite bottou and the other paper)

69 The stochastic gradient descent described above is a part of a much larger family of subgradient

70 based stochastic optimizers. One can, largely, avoid the issues related to choosing a good learning

71 rate by utilizing techniques that are robust to the choice of learning rate (cite adagrad and implicit-

72 implicit sgd). These stochastic subgradient algorithms will play a central role in parameter inference

73 for the four scenarios presented in table 1.

74 **shortcomings** In this work, we will focus on SGD algorithm for two major reasons. 1) The

75 simplicity of the algorithm and 2) its application to all four settings represented by table 2. However,

76 SGD is not always the appropriate optimization tool. Aside from it's slower rate of convergence

77 compared to second order methods, SGD is 1) reliant on choosing a good step-size and 2) is unable

78 to deal with sparsity inducing, $\ell_1$ regularization.

79 **rebuttal to shortcomings** Fortunately, SGD algorithm can be modified to remedy both of the
80 aforementioned shortcomings. The reliance on step-size can be reduced by modifying the algorithm
81 to adaptively learn a per-dimension step-size (5; 22; 9). To address the second shortcoming, various
82 algorithms and modifications to the vanilla SGD algorithm can be used to obtain a sparse solution
83 (6; 12; 20).

84 In what follows, we will explore each setting of table1. For each setting, we offer a short literature
85 review followed by a few relevant algorithms. While we cover all four settings, but we will focus our
86 attention on offline-distributed and online-distributed.

## Setting 1: Offline and Centralized

88 **Problem set up and solution** is the most common setting for statistical analysis by researchers.
89 We define this setting as the case where all the data is available at the time of computation (therefore
90 offline), and at a single computational node. For a GLM model, the (log)-likelihood can be efficiently
91 evaluating for a parameter vector $\hat{\beta}$. Therefore, for a small enough problem, the solution can be
92 found using Newton-Raphson method. However, since this method requires inverting the Hessian
93 matrix, it scales poorly for large number of covariates (suggest papers for that use). A closely
94 related optimization technique is the Fisher Scoring method which replaces the Hessian in Newton's
95 method with the expected value of Hessian. Using Fisher Scoring method, a model can be solved
96 by iteratively solving weighted linear regression problems (iteratively weighted least squares)

97 While the rate of convergence is higher for both of the aforementioned methods (cite) stochastic
98 optimization techniques still provide a plausible alternative particularly when the dataset is large.
99 We will further discuss these methods in the next 3 sections.

## Setting 2: Offline and Distributed

101 **related lit** In this setting, all the required data is gathered but the data is distributed across
102 multiple nodes and cannot be centralized. This setting naturally arises when the communication
103 of data is too costly or privacy is of concern. (8) provides an algorithm for distributed logistic
104 regression but considers the setting where the covariates are distributed across the nodes. (16)
105 provides three novel algorithms for lasso regression. In all three algorithms, the parameters are
106 computed locally and updated based on the parameters computed on "neighboring" datasets. In
107 general, Alternating Direction Method of Multipliers (ADMM) provides a natural framework for
108 consensus optimization of convex functions (1). In this method the number of parameters is first

expanded to include a set of private parameters for each node as well as a set of public parameters then the optimization problem is modified to optimize over the local parameters with the constraint that the local parameters must match the global parameters. Using this framework, the solution can be iteratively computed as detailed in (1; 14).

**Algorithm** Another promising approach to solving the problem is via a simple adaptation of the SGD algorithm presented earlier. Let $D_1 \ldots D_M$ denote the DO's then, Algorithm 1 represents a simple way of performing SGD without gathering the data. In any practical implementation, a few caveats need to be considered.

**mini-batch 1)** Due to network latency, it may be beneficial to compute the gradient on a batch size $> 1$ to reduce number of communications with the DO.

**parallelization 2)** Algorithm 1 assumes the information exchange happens sequentially. In practice, this algorithm can be parallelized. Unfortunately, the parallelization algorithms that assume the data can be assigned to nodes in a balanced way, may not always be applicable (23). The main concern with parallelization is that as one DO updates the parameter the other DO's may be busy computing the gradient using an older version of the parameter. One possibility is to use $k$ DO's, compute $\frac{b}{k}$ gradients in a parallel at each DO and synchronously combine the results into an average

125 gradient from a batch of size $b$ (3; 2).

---

**Algorithm 1:** SGD on distributed data

> **Result:** $\beta$
>
> $\beta_0, t \leftarrow 0, \eta_0, \text{converged} \leftarrow \text{FALSE}$
>
> **while** Not converged **do**
>> **for** $i$ $in$ $1 \dots M$ **do**
>>> **for** $x_j$ $in$ $D_i$ **do**
>>>> Send $\beta_t$ to $D_i$
>>>>
>>>> Retrieve $\nabla C_i(\beta_t)$ computed at $D_i$
>>>>
>>>> $\beta_{t+1} \leftarrow \beta_t - \eta_t \nabla C_i(\beta_t)$
>>>>
>>>> $t \leftarrow t + 1$
>>>>
>>>> **if** *condition* **then**
>>>>> converged = TRUE
>>>>>
>>>>> break
>>>>
>>>> **end**
>>>
>>> **end**
>>
>> **end**
>
> **end**

---

127 **Short comings and Extensions**   Stochastic gradient descent is readily expendable to online
128 learning

129 **SGD in distributed setting**

130 **Algorithm**   If parallelization is not an issue

131 **Setting 3,4: Online**

132 When data arrives one (or few) at a time, one may choose to rerun the entire regression in an offline
133 setting with the currently available data, however, this approach may be costly and impractical if
134 new data frequently becomes available. Online learning algorithms attempt to compute an efficient
135 update to the previous estimates using the information from the new data point. Bayesian frame-
136 work provides a natural way of updating the model parameters as new data arrives. In this setting
137 common practices for computing the update include, numerical integration as in (24) or replacing

138  the true posterior distribution with an approximate posterior chosen from a parametric distribution.

139  (19; 21; 18)

140      Another common approach is to use sub-gradient optimization methods to calculate online

141  updates for the model parameter(5; 12; 10). A simple algorithm, in this setting, would update

142  compute the gradient as a new point arrives, and update the parameter estimates based on this

143  gradient but, in contrast with algorithm 1, it would not loop over the data to reach convergence.


# Results

145  In this section, we will demonstrate the performance of linear regression and logistic regression

146  optimized using a selected set of algorithms from the list reviewed in the previous section. We

147  will generate synthetic data as explained in the next section to evaluate the performance of each

148  algorithm for varying number of covariates, data sizes, and covariance structures.


## DataSets

150  To benchmark the algorithms presented earlier, we will simulate dataset as follows:

151      Uncorrelated continuous covariates are drawn from a $Norm(0, 10^\alpha)$ with $\alpha$ drawn unifromly from

152  the interval $(0, 3)$. The discrete covariates are limited to take dosage values ($\{0,1,2\}$). The dosage

153  values are assigned by drawing a frequency parameter for each covariate ($f \sim Unif(0.01, 0.49)$) and

154  assigning a dosage of 0,1,2 with probabilities $f^2, 2f(1-f),$ and $(1-f)^2$ for each individual. Let $X_i$

155  denote the row vector of the covariates for individual $i$, and $\beta$ denote the vector of coefficients (drawn

156  from a $Norm(0, 1)$ distribution). Using these definitions, the label for individual $i$ is computed as

157  $y_i = X_i\beta + b + \epsilon_i$ for linear regression and $y_i = binom(1, X_i\beta + b + \epsilon_i)$ for logistic regression. In both

158  cases, $\epsilon_i$ is an individual-specific Gaussian error term with $Norm(0, 0.1)$ distribution and the bias

159  term ($b$) is set to zero ($b = 0$). The data is distributed across 5 DO's. The first four DO receive a

160  random number of individuals (distributed according to $Pois(\frac{\text{size}}{\#\text{DO's}})$) and the last DO received the

161  remaining individuals. The simulation is restarted if any DO receives less than 10% of the data.

162      We study the effects of correlation for continuous variables only. The covariates are drawn from

163  a multivariate normal random variable with a random covariance matrix generated via vine method

164  detained in (13). The remaining values are assigned as before.

## Algorithms

In order to benchmark the general approaches outlined in Methods, we implemented benchmarked 7 representative algorithms. Since we anticipate the largest cost to be the communication between the DO's and the central hub, we will evaluate these algorithms based on the number of communications rather than the total compute time. Below is a short list of the implemented algorithms:

**SGD** is the simple stochastic gradient algorithm introduced earlier. The stepsize is chosen to be constant for the first 500 iterations and will decays as $\eta_t = \frac{\eta}{(t-500+1)^{0.51}}$ for the remaining steps.

**ADAGRAD** uses the algorithm presented in (5). This algorithm uses previous gradient information to estimate the currect stepsize for each covariate. As a result, it is more robust to stepsize choices and can have a per-covariate stepsize.

**RMSPROP** uses the algorithm presented in (). This algorithm is similar to ADAGRAD but puts a heavier emphasis on the recent gradients in hopes of recucing the heavy reduction step-size for larger iterations.

**ADADELTA** uses the algorithm presented in (22). This algorithm also uses previous gradients and updates to produce a covariate-specifid step-size.

**SQN** uses the algorithm presented in () (byrd). In this algorithm every few iterations a noisy Hessian is estimated using a relatively large batch size. This Hessian is then used to compute the per-covariate step-size for each stochastic gradient descent step.

**AVG** computes the model parameters at each DO and uses a weighted average to compute the overall model parameter. As a result, this algorithm only has one iteration.

**ADMM** implements the ADMM algorithm (1). We will use an $\ell_1$ normalized logistic regression and LASSO () regression with very small penalties in lieu of logistic regression and oridinary least squares.

More information about the parameters, stepsizes, and implementations can be found in **reference link**.

## Performance

In order to benchmark these algorithms, we restricted the number of communications to 1,000 and measured the performance for different number of covariates drawn from purely gaussian, equal mix of guassian and dosage and purely dosage distributions using 5 DO's. For each experiment, the data was gathered to compute the gold-standard solution. Figure 3.4 shows the per-sample logistic cost
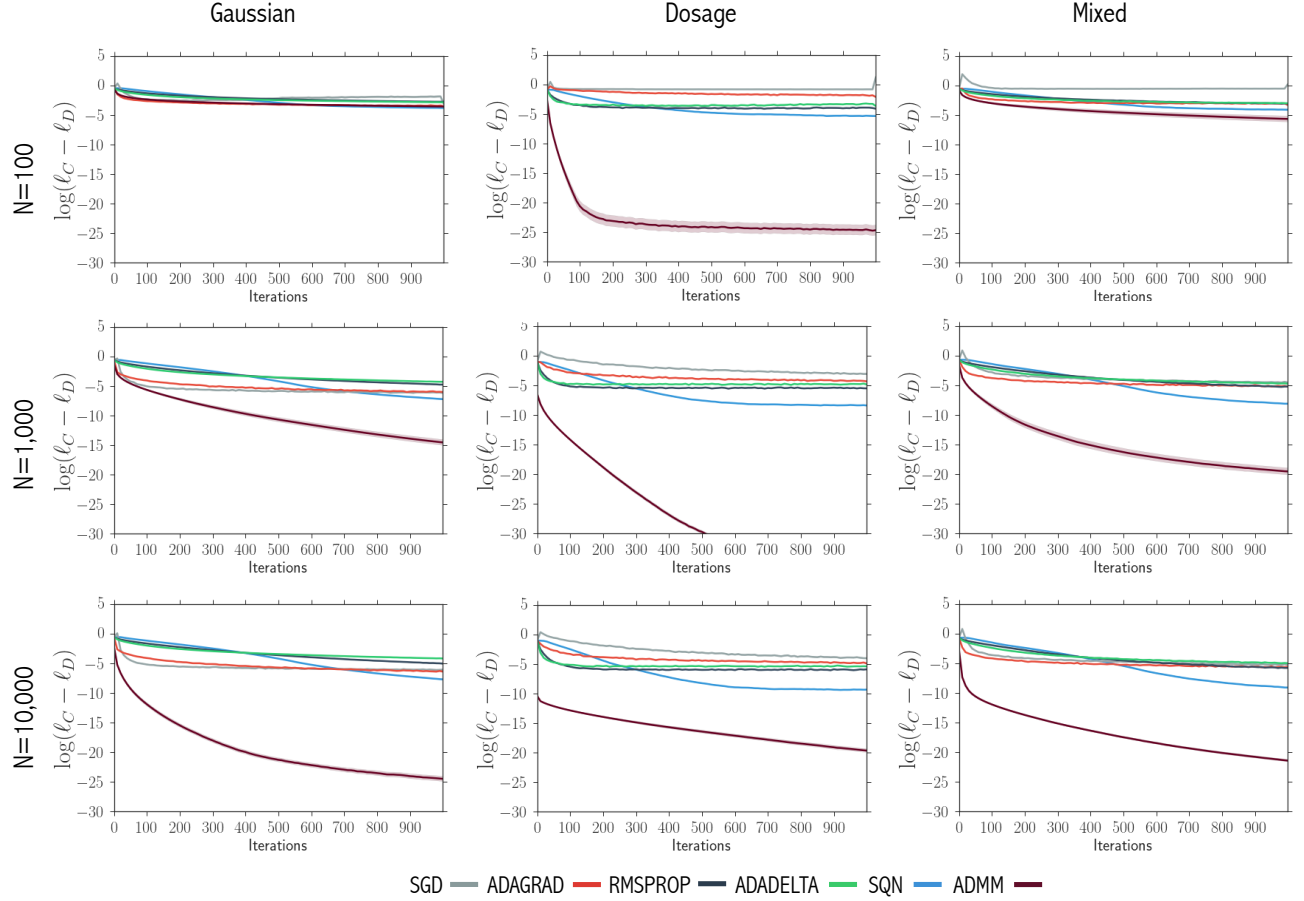
Figure 1: Plots of convergence of each method to the solution from a centralized logistic regression for 9 covariates (+1 intercept). In all cases, ADMM approaches the solution from centralized logistic regression faster than the other methods. Furthermore, at 1,000 iterations, the ADMM solution outperforms all other methods. As expected, ADMM converges much faster when the amount of data in each DO is large.

function difference for these algorithms and the combined logistic regression for 2,3,6,11 covariates (including the intercept value). Figure **??** shows a similar result for least squares regression and per-sample mean-squared error.

We note that ADMM and SQN consistently outperform other algorithms. Simple averaging performs well, when the sample size is large compared to the number of covariates, however, when this is not the case, it can produce very inaccurate estimates and therefore had to be removed from the top row of figure 3.4.
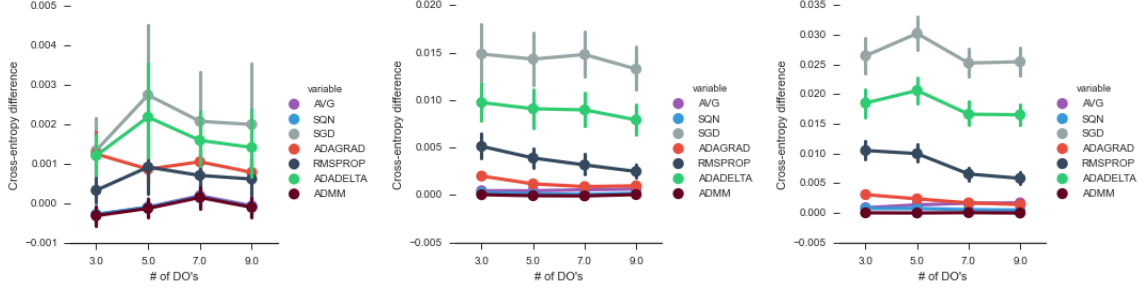
**PUT STUFF RELATED TO LEAST SQUARES HERE**

Figure 2: Insert caption w/ batch size and ...

## Performance as a function of number of DO's

Once again, we see that when the per-DO sample size is large compared to the number of covariates, simply averaging the per-DO estimates is a very accurate estimator, however this estimator fails when the per-DO estimates deteriorate. On the other hand SQN and ADMM are accurate throughout the parameter space explored.

## Performance as a function of number of communications

In the previous experiments, we have limited all the algorithms to 1000 iterations.

# Discussion

# Conclusions

# Figures and Tables

Figure 3: **This is the 1-line title of the figure.**

This is a longer text for detailed explanation. Should be longer than 1-line.

Table 3: **This is the 1-line title of the table.**

This is a longer text for detailed explanation. Should be longer than 1-line.

# References

1 Boyd S. 2011. Alternating direction method of multipliers. In: Talk at NIPS Workshop on Optimization and Machine Learning.

2 Chen J, Monga R, Bengio S, Jozefowicz R. 2016. Revisiting distributed synchronous sgd. arXiv preprint arXiv:1604.00981 .

3 Dekel O, Gilad-Bachrach R, Shamir O, Xiao L. 2012. Optimal distributed online prediction using mini-batches. Journal of Machine Learning Research 13:165–202.

4 Domingos P. 2012. A few useful things to know about machine learning. Communications of the ACM 55:78–87.

5 Duchi J, Hazan E, Singer Y. 2011. Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research 12:2121–2159.

6 Duchi J, Shalev-Shwartz S, Singer Y, Chandra T. 2008. Efficient projections onto the l 1-ball for learning in high dimensions. In: Proceedings of the 25th international conference on Machine learning. ACM. p. 272–279.

7 Friedman J, Hastie T, Tibshirani R. 2010. Regularization paths for generalized linear models via coordinate descent. Journal of statistical software 33:1.

8 Gopal S, Yang Y. 2013. Distributed training of large-scale logistic models. In: ICML (2). p. 289–297.

9 Kingma D, Ba J. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 .

10 Langford J, Li L, Strehl A. 2007. Vowpal wabbit online learning project.

11 Langford J, Li L, Strehl A. 2011. Vowpal wabbit. URL https://github.com/JohnLangford/vowpal wabbit/wiki .

12 Langford J, Li L, Zhang T. 2009. Sparse online learning via truncated gradient. Journal of Machine Learning Research 10:777–801.

238  13 Lewandowski D, Kurowicka D, Joe H. 2009. Generating random correlation matrices based on
239     vines and extended onion method. Journal of multivariate analysis 100:1989–2001.

240  14 Lubell-Doughtie P, Sondag J. 2013. Practical distributed classification using the alternating
241     direction method of multipliers algorithm. In: Big Data, 2013 IEEE International Conference
242     on. IEEE. p. 773–776.

243  15 Marschner IC, et al. 2011. glm2: fitting generalized linear models with convergence problems.
244     The R journal 3:12–15.

245  16 Mateos G, Bazerque JA, Giannakis GB. 2010. Distributed sparse linear regression. IEEE
246     Transactions on Signal Processing 58:5262–5276.

247  17 Nelder JA, Baker RJ. 1972. Generalized linear models. Encyclopedia of statistical sciences .

248  18 Opper M. 1996. On-line versus off-line learning from random examples: General results. Physical
249     Review Letters 77:4671.

250  19 Opper M, Winther O. 1999. A bayesian approach to on-line learning. On-line Learning in Neural
251     Networks, ed. D. Saad :363–378.

252  20 Shalev-Shwartz S, Tewari A. 2011. Stochastic methods for l1-regularized loss minimization.
253     Journal of Machine Learning Research 12:1865–1892.

254  21 Winther O, Solla SA. 1998. Optimal bayesian online learning. Theoretical Aspects of Neural
255     Computation (TANC-97), KYM Wong, I. King and D.-Y. Yeung eds., Springer Verlag, Singapore
256     .

257  22 Zeiler MD. 2012. Adadelta: an adaptive learning rate method. arXiv preprint arXiv:1212.5701 .

258  23 Zinkevich M, Weimer M, Li L, Smola AJ. 2010. Parallelized stochastic gradient descent. In:
259     Advances in neural information processing systems. p. 2595–2603.

260  24 Zoeter O. 2007. Bayesian generalized linear models in a terabyte world. In: Image and Signal
261     Processing and Analysis, 2007. ISPA 2007. 5th International Symposium on. IEEE. p. 435–440.