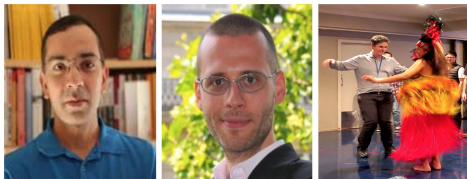


Preconditioning an Artificial Neural Network Using Naive Bayes

Nayyar A. Zaidi, Francois Petitjean, Geoffrey I. Webb



MONASH University
Information Technology

Introduction

- ▶ Maximum likelihood estimates of naive Bayes probabilities can be used to greatly speed-up logistic regression
- ▶ This talk demonstrates that this speed-up can also be attained for Artificial Neural Networks
- ▶ Talk outline
 - ▶ Introduction (2 minutes)
 - ▶ Proposed Approach (6 minutes)
 - ▶ Experimental Analysis (6 minutes)
 - ▶ Future Research, Q & A (3 minutes)

Contributions of the Paper

- ▶ We show that:
 1. Preconditioning based on naive Bayes is applicable and equally useful for Artificial Neural Networks (ANN) as it is for Logistic Regression (LR)
 2. Optimizing MSE objective function leads to lower bias than optimizing CLL, this leads to lower 0-1 Loss and RMSE on big datasets

Logistic Regression

- ▶ One of the state-of-the-art classifier
- ▶ Maximizes the Conditional Log-Likelihood (NLL)

$$\text{CLL}(\beta) = \sum_{i=1}^N \log P_{LR}(y^{(i)} | \mathbf{x}^{(i)}) \quad (1)$$

- ▶ If constrained to categorical attributes and multi-class problems, it leads to:

$$P_{LR}(y | \mathbf{x}) = \frac{\exp(\beta_y + \sum_{i=1}^a \beta_{y,i,x_i})}{\sum_{c \in \Omega_Y} \exp(\beta_c + \sum_{j=1}^a \beta_{c,j,x_j})} \quad (2)$$

- ▶ and

$$\exp\left(\beta_y + \sum_{i=1}^a \beta_{y,i,x_i} - \log \sum_{c \in \Omega_Y} \exp\left(\beta_c + \sum_{j=1}^a \beta_{c,j,x_j}\right)\right) \quad (3)$$

Naive Bayes and Weighted Naive Bayes

- ▶ Naive Bayes can be written as:

$$P_{NB}(y | \mathbf{x}) = \frac{P(y) \prod_{i=1}^a P(x_i | y)}{\sum_{c \in \Omega_Y} P(c) \prod_{j=1}^a P(x_j | c)} \quad (4)$$

- ▶ Adding weights in naive Bayes:

$$P_W(y | \mathbf{x}) = \frac{P(y)^{w_y} \prod_{i=1}^a P(x_i | y)^{w_{y,i,x_i}}}{\sum_{c \in \Omega_Y} P(c)^{w_c} \prod_{j=1}^a P(x_j | c)^{w_{c,j,x_j}}} \quad (5)$$

$$= \exp\left(w_y \log P(y) + \sum_{i=1}^a w_{y,i,x_i} \log P(x_i | y) - \right.$$

$$\left. \log \sum_{c \in \Omega_Y} \exp\left(w_c \log P(c) + \sum_{j=1}^a w_{c,j,x_j} \log P(x_j | c)\right)\right). \quad (6)$$

▶ LR

$$\exp\left(\beta_y + \sum_{i=1}^a \beta_{y,i,x_i} - \log \sum_{c \in \Omega_Y} \exp\left(\beta_c + \sum_{j=1}^a \beta_{c,j,x_j}\right)\right)$$

▶ Weighted Naive Bayes

$$\exp\left(w_y \log P(y) + \sum_{i=1}^a w_{y,i,x_i} \log P(x_i|y) - \log \sum_{c \in \Omega_Y} \exp\left(w_c \log P(c) + \sum_{j=1}^a w_{c,j,x_j} \log P(x_j|c)\right)\right)$$

▶ LR

$$\exp\left(\beta_y + \sum_{i=1}^a \beta_{y,i,x_i} - \log \sum_{c \in \Omega_Y} \exp\left(\beta_c + \sum_{j=1}^a \beta_{c,j,x_j}\right)\right)$$

▶ Weighted Naive Bayes

$$\exp\left(w_y \log P(y) + \sum_{i=1}^a w_{y,i,x_i} \log P(x_i | y) - \log \sum_{c \in \Omega_Y} \exp\left(w_c \log P(c) + \sum_{j=1}^a w_{c,j,x_j} \log P(x_j | c)\right)\right)$$

WANBIA-C

- ▶ LR

$$\exp\left(\beta_y + \sum_{i=1}^a \beta_{y,i,x_i} - \log \sum_{c \in \Omega_Y} \exp\left(\beta_c + \sum_{j=1}^a \beta_{c,j,x_j}\right)\right)$$

- ▶ Weighted Naive Bayes

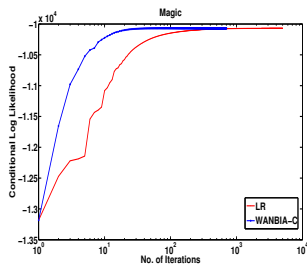
$$\exp\left(w_y \log P(y) + \sum_{i=1}^a w_{y,i,x_i} \log P(x_i | y) - \log \sum_{c \in \Omega_Y} \exp\left(w_c \log P(c) + \sum_{j=1}^a w_{c,j,x_j} \log P(x_j | c)\right)\right)$$

- ▶ $\beta_c \propto w_c \log P(c)$ and $\beta_{c,i,x_i} \propto w_{c,i,x_i} \log P(x_i | c)$
- ▶ WANBIA-C: Proposed in [1] shows an equivalence between LR and weighted naive Bayes
- ▶ For sake of clarity - we denote it by: WANBIA_{CLL}^C

WANBIA-C

- ▶ View 1: Learn weights by optimizing CLL to alleviate naive Bayes independence assumption
- ▶ View 2: $\text{WANBIA}_{\text{CLL}}^{\text{C}}$ uses generative estimates of the probabilities to speed-up the convergence
- ▶ View 3: Way of combining generative and discriminative models

- ▶ $\text{WANBIA}_{\text{CLL}}^{\text{C}}$ and LR generate equivalent models
- ▶ But have different convergence profiles



Artificial Neural Networks (LR)

- ▶ Minimizes MSE Objective Function instead of NLL
- ▶ We begin by writing an objective function:

$$\text{MSE}(\beta) = \frac{1}{2} \sum_{i=1}^N \sum_{c=1}^C \left(\delta(y = c) - P(c|\mathbf{x}^{(i)}) \right)^2$$

- ▶ where

$$P(c | \mathbf{x}) = \frac{\exp(\beta_c + \sum_{i=1}^a \beta_{c,i,x_i})}{\sum_{c' \in \Omega_Y} \exp(\beta_{c'} + \sum_{j=1}^a \beta_{c',j,x_j})}$$

Artificial Neural Networks (WANBIA-C)

- ▶ Minimizes MSE Objective Function instead of NLL
- ▶ We begin by writing an objective function:

$$\text{MSE}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \sum_{c=1}^C \left(\delta(y = c) - P(c|\mathbf{x}^{(i)}) \right)^2$$

- ▶ where

$$P(c | \mathbf{x}) = \frac{P(y)^{w_y} \prod_{i=1}^a P(x_i | y)^{w_{y,i,x_i}}}{\sum_{c \in \Omega_Y} P(c)^{w_c} \prod_{j=1}^a P(x_j | c)^{w_{c,j,x_j}}}$$

Proposed Method – WANBIA_{MSE}^C

▶ Step 1:

- ▶ Calculate class-probabilities as $P(y) = \pi_y = \frac{\#_y+m/C}{N+m}$
- ▶ Calculate other probabilities as $P(x_i | y) = \theta_{x_i|c} = \frac{\#_{x_i,y}+m/|x_i|}{\#_y+m}$

▶ Step 2:

- ▶ Optimize MSE based on weighted naive Bayes
- ▶ Use gradient-based iterative optimization algorithm
- ▶ Calculate the gradient:

$$\frac{\partial \text{MSE}(\mathbf{w})}{\partial w_{k,i,x_i}} = - \sum_{i=1}^N \sum_c^C (\delta(y=c) - P(c|\mathbf{x})) \frac{\partial P(c|\mathbf{x})}{\partial w_{k,i,x_i}},$$

- ▶ where

$$\frac{\partial P(c|\mathbf{x})}{\partial w_{k,i,x_i}} = P(c|\mathbf{x}) (\delta(c=k) - P(k|\mathbf{x})) \delta(x_i) \log \theta_{x_i|k},$$

- ▶ Use L-BFGS to get parameter vector w

WANBIA_{MSE}^C vs. ANN

- ▶ Gradient of parameters can be defined as:

$$\frac{\partial \text{MSE}(\mathbf{w})}{\partial w_{k,i,x_i}} = - \sum_{i=1}^N \sum_c^C \left(\delta(y=c) - \hat{P}(c|\mathbf{x}) \right) P(y|\mathbf{x}) \\ (\delta(y=k) - P(k|\mathbf{x})) \log \theta_{x_i|k} \delta(x_i)$$

- ▶ Note for ANN, we have:

$$\frac{\partial \text{MSE}(\mathbf{w})}{\partial \beta_{k,i,x_i}} = - \sum_{i=1}^N \sum_c^C (\delta(y=c) - P(c|\mathbf{x})) P(y|\mathbf{x}) \\ (\delta(y=k) - P(k|\mathbf{x})) \delta(x_i)$$

- ▶ WANBIA_{MSE}^C has the effect of re-scaling the gradients of ANN

$$\frac{\partial \text{MSE}(\mathbf{w})}{\partial w_{k,i,x_i}} = \frac{\partial \text{MSE}(\beta)}{\partial \beta_{k,i,x_i}} \log \theta_{x_i|k}, \quad \frac{\partial \text{MSE}(\mathbf{w})}{\partial w_k} = \frac{\partial \text{MSE}(\beta)}{\partial \beta_k} \log \pi_k$$

Experimental Results

- ▶ 73 standard UCI datasets
- ▶ Algorithms evaluated in terms of bias, variance, 0-1 Loss and RMSE
- ▶ 40 datasets with < 1000 instances
- ▶ 21 between 1000 and 10000 instances
- ▶ 12 datasets with > 10000 instances
- ▶ Datasets are divided into two categories *All* and *Big*
- ▶ MDL discretization is used to discretize numeric attributes
- ▶ L-BFGS solver is used

MSE vs. CLL

	WANBIA _{MSE} ^C vs. WANBIA _{CLL} ^C		ANN vs. LR	
	W-D-L	p	W-D-L	p
	<i>All Datasets</i>			
Bias	39/15/18	0.007	36/14/22	0.086
Variance	21/8/42	0.011	26/7/38	0.168
0-1 Loss	33/12/27	0.519	34/9/29	0.614
RMSE	30/5/37	0.463	28/4/40	0.181
	<i>Big Datasets</i>			
0-1 Loss	10/1/1	0.011	8/2/2	0.109
RMSE	8/1/3	0.226	8/0/4	0.387

Table: Win-Draw-Loss: WANBIA_{MSE}^C vs. WANBIA_{CLL}^C and ANN vs. LR. p is two-tail binomial sign test. Results are significant if $p \leq 0.05$.

WANBIA_{MSE}^C vs. ANN

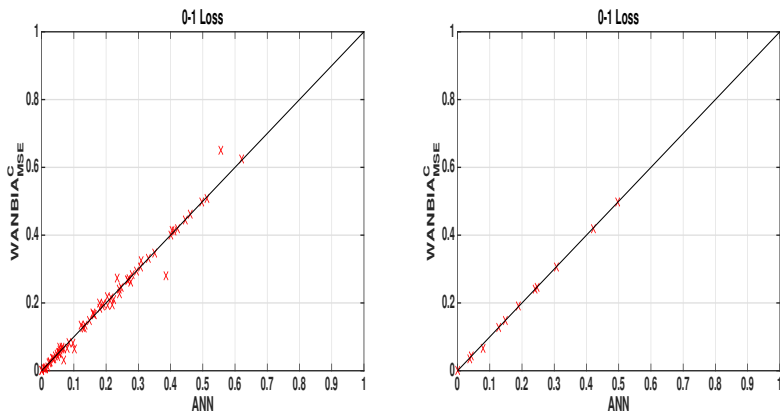


Figure: Comparative scatter of 0-1 Loss of ANN and WANBIA_{MSE}^C on *All* (Left) and *Big* (Right) datasets.

WANBIA_{MSE}^C vs. ANN

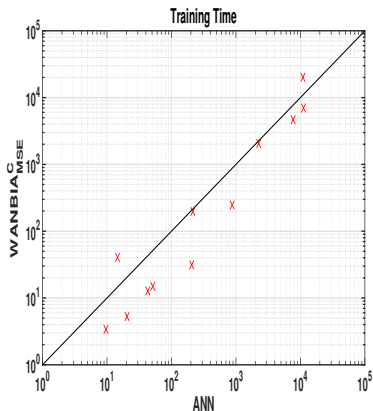
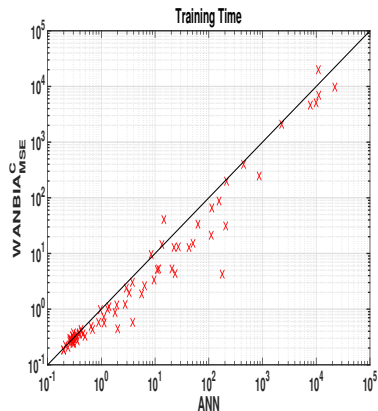


Figure: Comparative scatter of training time of ANN and WANBIA_{MSE}^C on *All* (Left) and *Big* (Right) datasets.

WANBIA_{MSE}^C vs. ANN

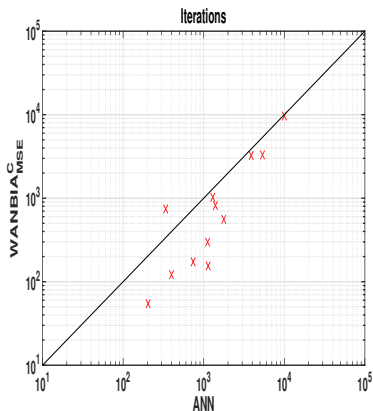
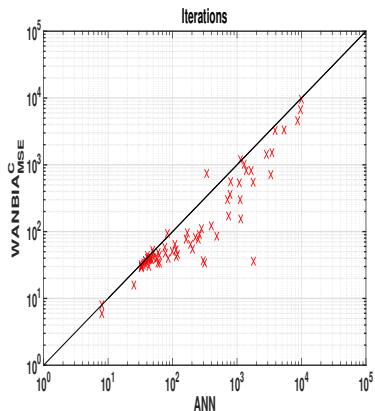


Figure: Comparative scatter of number of iterations to convergence of ANN and WANBIA_{MSE}^C on *All* (Left) and *Big* (Right) datasets.

Convergence Curves

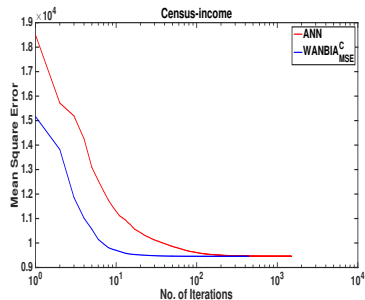
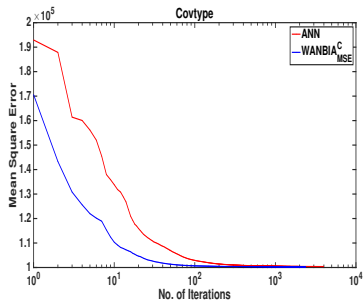


Figure: Comparative convergence profiles of ANN and WANBIA^C_{MSE} on Covtype (Left) and Census-income (Right) datasets.

Convergence Curves

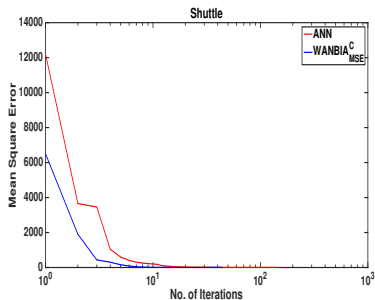
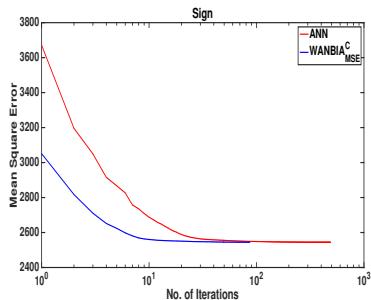
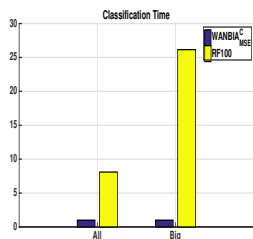
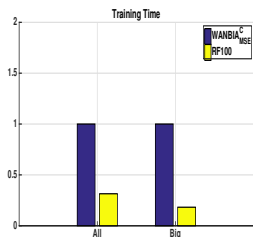


Figure: Comparative convergence profiles of ANN and WANBIA_{MSE}^C on Sign (Left) and Shuttle (Right) datasets.

WANBIA_{MSE}^C vs. Random Forest



WANBIA _{MSE} ^C vs. RF100		
	W-D-L	p
<i>All Datasets</i>		
Bias	41/5/26	0.086
Variance	32/2/38	0.550
0-1 Loss	30/2/40	0.282
RMSE	27/0/45	0.044
<i>Big Datasets</i>		
0-1 Loss	5/0/7	0.774
RMSE	5/0/7	0.774





Table: Win-Draw-Loss: WANBIA_{MSE}^C vs. Random Forest. p is two-tail binomial sign test. Results are significant if $p \leq 0.05$.

Conclusion and Future Work

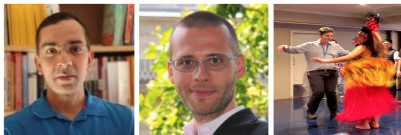
- ▶ Simple (naive Bayes based) preconditioning can speed-up convergence of ANN
- ▶ The proposed $\text{WANBIA}_{\text{MSE}}^{\text{C}}$ approach has the desirable property of asymptoting to optimum much quicker than ANN
- ▶ We are investigating:
 1. Why naive Bayes estimates are such a good pre-conditioner?
 2. An out-of-core Stochastic Gradient Descent (SGD) optimization
 3. $\text{WANBIA}_{\text{MSE}}^{\text{C}}$ for ANN with hidden layers
 4. Applicability of WANBIA-C style pre-conditioning to other objective functions

▶ Q & A

▶ Offline Discussions

- ▶  @nayyar_zaidi
- ▶  nayyar.zaidi@monash.edu
- ▶  nayyar_zaidi
- ▶  <http://users.monash.edu.au/~nzaidi>

▶ For further discussions, contact:



MONASH University
Information Technology



N. A. Zaidi, M. J. Carman, J. Cerquides, and G. I. Webb, “Naive-Bayes inspired effective pre-conditioners for speeding-up logistic regression,” in *IEEE International Conference on Data Mining*, pp. 1097–1102, 2014.



N. A. Zaidi, J. Cerquides, M. J. Carman, and G. I. Webb, “Alleviating naive Bayes attribute independence assumption by attribute weighting,” *Journal of Machine Learning Research*, vol. 14, pp. 1947–1988, 2013.