120??

Theorem [theorem]Corollary [theorem]Definition [theorem]Lemma [theorem]Proposition [the-

rem]Remark

The Canadian Journal of Statistics **Y**ol. xx, No. yy, 20??, Pages 1–25 La revue canadienne de statistique Weighted **Bayesian Boot**strap for Scalable **Bayes** BlindedA1* BlindedB² $^{1}Author$ affiliаtions will go © 20?? Statistical Society of Canada / Société statistique du Canada

CJS ???

BLINDEDA AND BLANDEDA here in the accepted manuscript, but do NOTinclude them in your initial submission because it

must

520??	
be	
anonymous.	
² Second	
Affiliation	
Bayesian,	
Boot-	
strap,	
MCMC,	
Weighted	
Boot-	
strap,	
ABC,	
Trend	
Fil-	
ter-	
ing,	
Deep	
Learn-	
ing,	
Ten-	
sor-	
Flow,	
Reg-	
u-	

BLINDEDA ANIDIBŁKŅĪDEDJĄ larization: MSC 2010: Primary 62???; secondary 62??? We develop a weighted Bayesian Bootstrap (WBB) for machine

720??	
learn-	
ing	
and	
statis-	
tics.	
WBB	
pro-	
vides	
un-	
cer-	
tainty	
quan-	
tifi-	
са-	
tion	
by	
sam-	
pling	
from	
a	
high	
di-	
men-	
sional	
pos-	

BLINDEDA ANIDIBLIKŅ DEDA terior distribution. WBB is computationally fast and scalable using only offtheshelf opti-

920??	
miza-	
tion	
soft-	
ware	
such	
as	
Ten-	
sor-	
Flow.	
We	
pro-	
vide	
reg-	
u-	
lar-	
ity	
con-	
di-	
tions	
which	
ар-	
ply	
to	
a	
wide	

BOLINDEDA ANADIBŁKŅ DEDA range of machine learning and statistical models. We illustrate our methodology in regular-

1	20??	
	ized	
	re-	
	gres-	
	sion,	
	trend	
	fil-	
	ter-	
	ing	
	and	
	deep	
	learn-	
	ing.	
	Fi-	
	nally,	
	we	
	con-	
	clude	
	with	
	di-	
	rec-	
	tions	
	for	
	fu-	
	ture	
	re-	

B2:INDEDA ANIOIBLANIOEDA search. The Canadian Journal of Statistics xx: 1-25; 20?? © 20?? Statistical Society of Canada Résumé: Insérer votre résumé ici. We will

1320?? supply a French abstract for those authors who can't prepare it themselves. La revue canadienne de statistique xx: 1–

BALINDEDA ANADIBŁANDEDA

25; 20??

© 20?? Société

statis-

tique

du

Canada

*Author

to

whom

cor-

re-

spon-dence

may

be

addressed.

E-

mail:

In-

sert

your email

ad-

dress

here

only af-ter

your

ра-

per has

been ac-

cepted

1:	520??
	1.
	IN-
	TRO-
	DUC-
	TION
	Weighted
	Bayesian
	Boot-
	strap
	(WBB)
	is
	a
	simulation-
	based
	al-
	go-
	rithm
	for
	as-
	sess-
	ing
	un-

B6INDEDA AND BLKN 106 D/B certainty in machine learning and statistics. Uncertainty quantification (UQ) is an active

1	720??
	area
	of
	re-
	search,
	par-
	tic-
	u-
	larly
	in
	high-
	dimensional
	in-
	fer-
	ence
	prob-
	lems.
	Whilst
	there
	are
	com-
	pu-
	ta-

B&INDEDA ANADIBŁKŅ DEDA tionally fast and scalable algorithms for training models in a wide variety of

1	920??
	con-
	texts,
	un-
	cer-
	tainty
	as-
	sess-
	ment
	is
	still
	re-
	quired.
	De-
	vel-
	op-
	ing
	com-
	pu-
	ta-
	tion-
	ally
	fast

BOLINDEDA ANADIBŁKŅ DEDA scalable algorithms for sampling a posterior distribution is a notori-

2	20??
	ously
	hard
	prob-
	lem.
	WBB
	makes
	a
	con-
	tri-
	bu-
	tion
	to
	this
	lit-
	er-
	a-
	ture
	by
	show-
	ing
	how

BILINDEDA ANKO IBLANDEDA
off-
the-
shelf
op-
ti-
miza-
tion
al-
go-
rithms,
such
as
con-
vex
op-
ti-
miza-
tion
or
stochas-
tic

2	320??	
	gra-	
	di-	
	ent	
	de-	
	scent	
	(SGD)	
	in	
	Ten-	
	sor-	
	Flow	
	can	
	also	
	be	
	used	
	to	
	pro-	
	vide	
	un-	
	cer-	
	tainty	
	quan-	
	tifi-	

B4LINDEDA ANADIBŁKŅ NOEDA cation. Our work builds on ? who provide a weighted likelihood Bootstrap (WLB) method for Bayesian

2	520??
	in-
	fer-
	ence.
	They
	de-
	velop
	a
	weighted
	like-
	li-
	hood
	Boot-
	strap
	al-
	go-
	rithm
	to-
	gether
	with
	the
	ар-

BEINDEDA AND BLANDEDA propriate asymptotic analуsis to show that such an algorithm provides efficient pos-

2	720??
	te-
	rior
	sam-
	ples.
	Their
	boot-
	strap
	pro-
	ce-
	dure
	ex-
	ploits
	the
	fact
	that
	the
	pos-
	te-
	rior
	dis-
	tri-

B&INDEDA ANIQIBLKIN DEDJIG		
bu-		
tion		
cen-		
tered		
at		
the		
max-		
i-		
mum		
like-		
li-		
hood		
es-		
ti-		
mate		
(MLE)		
has		
a		
sec-		
ond		
or-		

2920??			
	der		
	ex-		
	pan-		
	sion		
	that		
	also		
	de-		
	pends		
	on		
	the		
	prior		
	and		
	its		
	deriva-		
	tive.		
	The		
	weighted		
	Bayesian		
	Boot-		
	strap		
	(WBB)		

BOLINDEDA ANDOIBŁKŅIDEDA		
cal-		
cu-		
lates		
a		
se-		
ries		
of		
pos-		
te-		
rior		
modes		
rather		
than		
MLEs.		
This		
has		
the		
ad-		
van-		
tage		
that		

3120??			
	high		
	di-		
	men-		
	sional		
	pos-		
	te-		
	rior		
	modes		
	are		
	read-		
	ily		
	avail-		
	able		
	par-		
	tic-		
	u-		
	larly		
	us-		
	ing		
	the		
	reg-		

82 INDEDA AND IBLANDEDA ularized estimates are fast to compute from convex optimization methods or stochas-

3320??			
	tic		
	gra-		
	di-		
	ent		
	de-		
	scent		
	(SGD)		
	for		
	neu-		
	ral		
	net-		
	work		
	ar-		
	chi-		
	tec-		
	tures		
	such		
	as		
	deep		
	learn-		
	ing.		

84LINDEDA ANADIBŁKŅ NIEDA By linking WLB and WBB, with modernday optimization to calibrate estimate, we

3520??			
	pro-		
	vide		
	a		
	frame-		
	work		
	for		
	un-		
	cer-		
	tainty		
	quan-		
	tifi-		
	ca-		
	tion.		
	Uncertainty		
	es-		
	ti-		
	mates		
	are		
	pro-		
	vided		
	at		

BIGINDEDA ANVOIBLANNO EDA little to no extra cost. Quantifying uncertainty is typically unavailable

3	720??
	in
	a
	purely
	reg-
	u-
	lar-
	iza-
	tion
	op-
	ti-
	miza-
	tion
	method.
	An-
	other
	fea-
	ture
	that
	is
	straight-
	for-

88. INDEDA ANADIBŁKŅ NO E. DA ward to add is a regularization path across hyperparameters. This is so much easier than

3920??			
	tra-		
	di-		
	tional		
	Bayesian		
	to		
	do		
	prior		
	sen-		
	si-		
	tiv-		
	ity		
	anal-		
	y-		
	sis		
	where		
	hyper-		
	parameters		
	are		
	hard		
	to		
	as-		

BOLINDEDA ANADIBŁKŅ DEDA sess. Rather we use predictive crossvalidation techniques. The rest of the paper is outlined as

4	120??
	fol-
	lows.
	Sec-
	tion
	2
	de-
	vel-
	ops
	our
	weighted
	Bayesian
	Boot-
	strap
	(WBB)
	al-
	go-
	rithm.
	Sec-
	tion
	3
	pro-

BAZINDEDA ANADIBIAN, NO E.D/B vides an application to high dimensional sparse regression, trend filtering and deep

4	320??
	learn-
	ing.
	WBB
	can
	also
	be
	ар-
	plied
	to
	Bayesian
	tree
	mod-
	els
	(?).
	Fi-
	nally,
	Sec-
	tion
	4
	con-
	cludes

BALINDEDA ANADIBŁKŅ DEDA with directions for future research. Areas for future study include Bootstrap filters

4520??			
	in		
	state-		
	space		
	mod-		
	els		
	(?)		
	and		
	com-		
	par-		
	i-		
	son		
	with		
	the		
	resampling-		
	sampling		
	per-		
	spec-		
	tive		
	to		
	se-		
	quen-		

BIGINDEDA ANADIBŁKŅ DEDA tial Bayesian inference (?), etc. 2. WEIGHTED **BAYESIAN** BOOT-**STRAP** Let yan nvector of outcomes, θ

4	720??
	de-
	notes
	a
	d-
	dimensional
	pa-
	ram-
	e-
	ter
	of
	in-
	ter-
	est
	and
	A
	a
	fixed
	$n \times$
	d
	ma-
	trix

BBLINDEDA ANVOIBLAN, NO EDA whose rows are the design points (or "features") a_i^T where we index observations by i

4	920??
	and
	pa-
	ram-
	e-
	ters
	by
	j.
	A
	large
	num-
	ber
	of
	ma-
	chine
	learn-
	ing
	and
	sta-
	tis-
	ti-
	cal

BOLINDEDA ANVOIBLANNI DE DAR problems can be expressed in the form where $l(y|\theta) =$ $\sum_{i=1}^{n} \log f(y_i; a_i^{\top} \theta)$ is a measure of

fit

(or

5120??		
	"em-	
	pir-	
	i-	
	cal	
	risk	
	func-	
	tion")	
	de-	
	pend-	
	ing	
	im-	
	plic-	
	itly	
	on	
	A	
	and	
	y.	
	The	
	penalty	
	func-	
	tion	

B2:INDEDA ANADIBLANDEDA or regularization term, $\lambda\phi(\theta)$, effects a favorable biasvariance tradeoff. We allow for

5320??			
	the		
	pos-		
	si-		
	bil-		
	ity		
	that		
	$\phi(heta)$		
	may		
	have		
	points		
	in		
	its		
	do-		
	main		
	where		
	it		
	fails		
	to		
	be		
	dif-		
	fer-		

BALINDEDA ANDOIBLANDEDA
en-
tiable.
Suppose
that
we
ob-
serve
data,
y =
(y_1,\ldots,y_n)
from
a
model
pa-
ram-
e-
ter-
ized
by
θ .
For

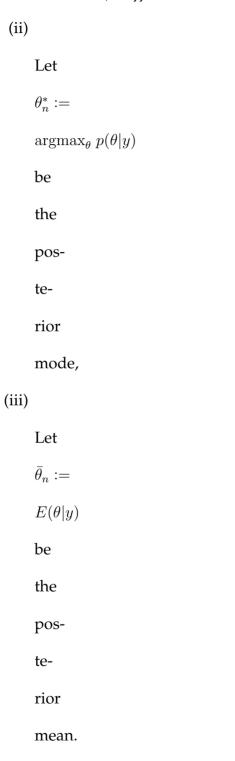
5520??			
	ex-		
	am-		
	ple,		
	we		
	might		
	have		
	a		
	prob-		
	a-		
	bilis-		
	tic		
	model		
	that		
	de-		
	pends		
	on		
	a		
	pa-		
	ram-		
	e-		
	ter,		
	θ ,		

BIGINDEDA ANIZOIBIAN, NIEDA where $p(y|\theta)$ is known as the likelihood function. Equivalently, we can define a measure

5720?? of fit $l(y|\theta) =$ $\log f(y;\theta) =$ $\log p(y|\theta).$ We will make use of the following (i) Let $\hat{\theta}_n :=$ $\operatorname{argmax}_{\theta} p(y|\theta)$ be the

MLE,

B&INDEDA ANIOIBLANDEDA



5920??			
	We		
	now		
	de-		
	velop		
	a		
	key		
	du-		
	al-		
	ity		
	be-		
	tween		
	reg-		
	u-		
	lar-		
	iza-		
	tion		
	and		
	pos-		
	te-		
	rior		
	boot-		

BOLINDEDA ANADIBŁKŅ DEDA strap simulation. 2.1. Bayesian Regularization Duality From the Bayesian perspective, the

6120?? measure of fit, $l(y|\theta) =$ $-\log f(y;\theta)$, and the penalty function, $\lambda\phi(\theta)$, correspond to the negative loga-

62 INDEDA AND IBLANDEDA rithms of the likelihood and prior distribution in the hierarchical

model

$$\begin{split} f(y;\theta) &= p(y|\theta) \propto \exp\{-l(y|\theta)\} \;, \quad p(\theta) \propto \exp\{-\lambda\phi(\theta)\} \end{split}$$

$$p(\theta|y) \propto \exp\{-(l(y|\theta) + \lambda\phi(\theta))\}. \end{split}$$

The

prior

is

not

nec-

es-

sar-

ily

proper

but

the

pos-

te-

rior,

 $p(\theta|y) \propto$

 $p(y|\theta)p(\theta)$,

may

BALINDEDA ANADIBLANDEDAB
still
be
proper.
This
pro-
vides
an
equiv-
a-
lence
be-
tween
reg-
u-
lar-
iza-
tion
and
Bayesian
meth-
ods.

6520??			
	For		
	ex-		
	am-		
	ple,		
	re-		
	gres-		
	sion		
	with		
	a		
	least		
	squares		
	log-		
	likelihood		
	sub-		
	ject		
	to		
	a		
	penalty		
	such		
	as		
	an		

66 INDEDA ANIOIBLAN DE DA L^2 norm (ridge) Gaussian probability model or L^1 norm (lasso) double exponential prob-

6720??

a-

bil-

ity

model.

We

then

have

$$\hat{\theta}_n = \lim_{\theta \in \Theta} l(y|\theta),$$
 (2)

$$\theta_n^* = \mathop{}_{\theta \in \Theta} \left\{ l(y|\theta) + \lambda \phi(\theta) \right\}$$

Let

 ∂

be

the

sub-

d-

if-

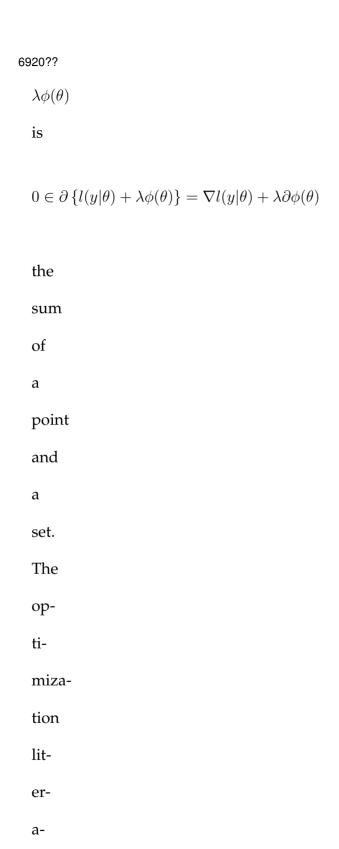
fer-

en-

tial

op-

68 INDEDA AND IBLAN DE DA erator. Then a necessary and sufficient condition for θ^* to minimize $l(y|\theta) +$



ture

BOLINDEDA ANVOIBLANDEDA

characterizes θ^* as the fixed point of a proximal operator $\theta^* =$ $\mathrm{prox}_{\gamma\phi}\{\theta^* -$

 $\lambda \nabla f(\theta^*)\},$

7120??			
	see		
	?		
	and		
	?		
	for		
	fur-		
	ther		
	dis-		
	cus-		
	sion.		
	A		
	gen-		
	eral		
	class		
	of		
	nat-		
	u-		
	ral		
	ex-		
	po-		
	nen-		

BAZINDEDA ANVOIBLANNO EDAS tial family models can be expressed in terms of the Bregman divergence of the dual

7	320??	
	of	
	the	
	cu-	
	mu-	
	lant	
	trans-	
	form.	
	Let	
	ϕ	
	be	
	the	
	con-	
	ju-	
	gate	
	Leg-	
	en-	
	dre	
	trans-	
	form	
	of	
	als	

BALINDEDA ANVOIBLANDEDA

_	_				
Н	4	2	n	C	ρ

$$\psi(\theta) =$$

$$\sup_{\mu} (\mu^{\top} \theta - \phi(\mu)).$$

Then

we

can

write

$$p_{\psi}(y|\theta) = \exp\left(y^{\top}\theta - \psi(\theta) - h_{\psi}(y)\right)$$
$$= \exp\left\{\inf_{\mu} \left((y - \mu)^{\top}\theta - \phi(\mu)\right) - h_{\psi}(y)\right\}$$
$$= \exp\left(-D_{\phi}(y, \mu(\theta)) - h_{\phi}(y)\right)$$

where

the

in-

fi-

mum

is

at-

tained

at

7!	7520??	
	$\mu(heta) =$	
	$\mu(\theta) = \phi'(\theta)$	
	is	
	the	
	mean	
	of	
	the	
	ex-	
	po-	
	nen-	
	tial	
	fam-	
	ily	
	dis-	
	tri-	
	bu-	
	tion.	
	We	
	rewrite	

 $h_{\psi}(y)$

is

BIGINDEDA ANADIBŁKŅ DEDA terms of the correction term and $h_{\phi}(y)$. Here there is a duality as D_{ϕ} can be in-

7	720??
	ter-
	preted
	as
	a
	Breg-
	man
	di-
	ver-
	gence.
	For
	a
	wide
	range
	of
	non-
	smooth
	ob-
	jec-
	tive
	func-

B&INDEDA ANIDIBŁKŅIDEDA tions/statistical models, recent regularization methods provide fast, scalable algorithms for

7	920??
	cal-
	cu-
	lat-
	ing
	es-
	ti-
	mates
	of
	the
	form
	(??),
	which
	can
	also
	be
	viewed
	as
	the
	pos-
	te-
	rior

BOLINDEDA ANVOIBLANNI DE DAR mode. Therefore as λ varies we obtain a full regularization path as a form of

8120??				
	prior			
	sen-			
	si-			
	tiv-			
	ity			
	anal-			
	y-			
	sis.			
	?			
	and			
	?			
	con-			
	sid-			
	ered			
	sce-			
	nar-			
	ios			
	where			
	pos-			
	te-			

BIZINDEDA ANIOIBILININI DEDIR
rior
modes
can
be
used
as
pos-
te-
rior
means
from
aug-
mented
prob-
a-
bil-
ity
mod-
els.
More-
over,

8	320??
	in
	their
	orig-
	i-
	nal
	foun-
	da-
	tion
	of
	the
	Weighted
	Like-
	li-
	hood
	Boot-
	strap
	(WLB),
	?
	in-
	tro-
	duced

B4LINDEDA ANIA IBŁKŅIOEDA
the
con-
cept
of
the
im-
plicit
prior.
Clearly
this
is
an
av-
enue
for
fu-
ture
re-
search.

8520??				
	2.2.			
	WBB			
	Al-			
	go-			
	rithm			
	We			
	now			
	de-			
	fine			
	the			
	weighted			
	Bayesian			
	Boot-			
	strap			
	(WBB).			
	Fol-			
	low-			
	ing			
	?,			
	we			
	con-			
	struct			

86 INDEDA ANVOIBLAN DIEDA

a

ran-

domly

weighted

pos-

te-

rior

dis-

tri-

bu-

tion

de-

noted

by

$$\mathbf{w} = (w_1, ..., w_n, w_p), \ p_{\mathbf{w}}(\theta|y) \propto \prod_{i=1}^n p(y_i|\theta)^{w_i} p(\theta)^{w_p}$$

where

the

weights

 $w_p, w_i \sim$

Exp(1)

8	720??
	are
	ran-
	domly
	gen-
	er-
	ated
	weights.
	It's
	equiv-
	a-
	lent
	to
	draw
	$w_i =$
	$\log(1/U_i)$
	where
	U_i 's
	are
	i.i.d.
	Uni-
	form
	(0,1),

88. INDEDA ANIOIBIAN, DEDA which is motivated by the uniform Dirichlet distribution for multinomial data. We

8920??				
	have			
	used			
	the			
	fact			
	that			
	for			
	i.i.d.			
	ob-			
	ser-			
	va-			
	tions,			
	the			
	like-			
	li-			
	hood			
	can			
	be			
	fac-			
	tor-			
	ized			
	as			

BOINDEDA ANTOIBERNINGEDAS $p(y|\theta) =$

 $\prod_{i=1}^n p(y_i|\theta).$

This

is

not

cru-

cial

for

our

anal-

у-

sis

but

is

a

com-

mon

as-

sump-

tion.

Let

9	20??	
	$ heta_{\mathbf{w},n}^*$	
	de-	
	note	
	the	
	mode	
	of	
	this	
	reg-	
	u-	
	lar-	
	ized	
	dis-	
	tri-	
	bu-	
	tion.	
	Again,	
	there	
	is	
	an	
	equiv-	
	a-	

82 INDEDA ANTO IBŁKŅ DEDA

lence

$$\theta_{\mathbf{w},n}^* := {}_{\theta} p_{\mathbf{w}}(\theta|y) \equiv {}_{\theta} \sum_{i=1}^{n} w_i l_i(y_i|\theta) + \lambda w_p \phi(\theta)$$

where

$$l_i(y_i|\theta) =$$

$$-\log p(y_i|\theta)$$

and

$$\lambda \phi(\theta) =$$

$$-\log p(\theta)$$
.

Note

that

we

have

a

weighted

like-

li-

hood

and

a

9320??		
	new	
	reg-	
	u-	
	lar-	
	iza-	
	tion	
	pa-	
	ram-	
	e-	
	ter,	
	λw_p .	
	The	
	crux	
	of	
	our	
	pro-	
	ce-	
	dure	
	is	
	to	
	cre-	

94LINDEDA ANADIBŁKŅ 1016ED/B ate a sample of the weighted posterior modes $\{\theta_{\mathbf{w},n}^*\}$ (computationally cheap as each sub-

9520??			
problem			
can			
be			
solved			
via			
op-			
ti-			
miza-			
tion).			
Our			
main			
re-			
sult			
is			
the			
fol-			
low-			
ing:			
Algorithm	ı :		
Weighted			

86 INDEDA ANIO IBLANDEDA Bayesian **Boot**strap (WBB) 1. Iterate: sample $\mathbf{w} =$ $\{w_1, w_2, ..., w_n, w_p\}$ via exponentials. $w_p, w_i \sim$ Exp(1). 2. For each

w,

solve

 $\theta_{\mathbf{w},n}^* =$

 $_{\theta}\sum_{i=1}^{n}w_{i}l_{i}(\theta)+$

 $\lambda w_p \phi(\theta)$.

The

WBB

al-

go-

rithm

is

fast

and

scal-

able

to

com-

pute

a

reg-

u-

lar-

98. INDEDA ANIO IBILIANIO ED/B ized estimator. For a large number of popular priors, the minimizing

9920??		
	SO-	
	lu-	
	tion	
	$ heta_{\mathbf{w},n}^*$	
	in	
	the	
	sec-	
	ond	
	step	
	can	
	be	
	di-	
	rectly	
	ob-	
	tained	
	via	
	reg-	
	u-	
	lar-	
	iza-	
	tion	
	pack-	

BOONDEDA AND IBLANNI EDA ages such as glmnet by Trevor Hastie and genlasso by Taylor Arnold. When the likelihood function or the

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

10	120??
	prior
	is
	spe-
	cially
	de-
	signed,
	Stochas-
	tic
	Gra-
	di-
	ent
	De-
	scent
	(SGD)
	is
	pow-
	er-
	ful
	and
	fast
	enough

BO2NDEDA ANIO IBŁKŅ DEDA to solve the minimization problem. It can be easily implemented in Tensor-Flow

10	0320??	
	once	
	the	
	ob-	
	jec-	
	tive	
	func-	
	tion	
	is	
	spec-	
	i-	
	fied.	
	See	
	Ap-	
	pendix	
	(??)	
	and	
	?	
	for	
	fur-	
	ther	
	dis-	
	cus-	

BO4NDEDA AND BLXNDEDB sion. The next section builds on ? and derives asymptotic properties of the weighted Bayesian

10	520??
	Boot-
	strap.
	We
	sim-
	ply
	add
	the
	reg-
	u-
	lar-
	ized
	fac-
	tor.
	То
	choose
	the
	amount
	of
	reg-
	u-
	lar-

BOSNDEDA AND BLANDEDA ization λ , we can use the marginal likelihood $m_{\lambda}(y)$, estimated by bridge sampling (?) or

10720??	
	sim-
	ply
	us-
	ing
	pre-
	dic-
	tive
	cross-
	validation.
	2.3.
	WBB
	Prop-
	er-
	ties
	The
	fol-
	low-
	ing
	propo-
	si-
	tion

BOSNDEDA ANTOIBLANDEDA which follows from the Theorem 2 in ? summaries the properties of WBB.

Proposition

Tioposition
The
weighted
Bayesian
Boot-
strap
draws
are
ар-
prox-
i-
mate
pos-
te-
rior
sam-
ples

 $\left\{\theta_{\mathbf{w},n}^{*(k)}\right\}_{k=1}^K \sim p(\theta|y).$

BLONDEDA ANVOIBLANDEDA Now we consider 'large n'properties. The variation in the posterior density $p(\theta|y) \propto$

11120?? $e^{-nl_n(\theta)}p(\theta)$ for sufficiently large nwill be dominated by the likelihood term. Expand-

ing

BL2NDEDA ANADIBŁKŅ DEDA $l_n(\theta)$ around its maximum, $\hat{\theta}$, and defining $J_n(\hat{\theta}) =$ $nj(\hat{\theta})$ as the observed infor-

ma-

tion

ma-

113	1320??			
	trix			
	gives			
	the			
	tra-			
	di-			
	tional			
	nor-			
	mal			
	ар-			
	prox-			
	i-			
	ma-			
	tion			
	for			
	the			
	pos-			
	te-			
	rior			
	dis-			
	tri-			
	bu-			

BL4NDEDA AND BLXNDEDB

tion

$$\theta \sim N_d \left(\hat{\theta}_n, J_n^{-1}(\hat{\theta}) \right)$$

where

 $\hat{\theta}_n$

is

the

MLE.

A

more

ac-

cu-

rate

ap-

prox-

i-

ma-

tion

is

ob-

tained

11	520??
	by
	ex-
	pand-
	ing
	around
	the
	pos-
	te-
	rior
	mode,
	θ^* ,
	which
	we
	will
	ex-
	ploit
	in
	our
	weighted
	Bayesian
	Boot-

BLENDEDA AND BLANDEDA strap. Now we have the asymptotic distributional approximation $\theta \sim N_d \left(\theta^*, J_n^{-1}(\theta^*) \right)$ where

 $\theta_n^* \vcentcolon=$

 $\arg \max_{\theta} p(\theta|y)$

117	1720??		
	is		
	the		
	pos-		
	te-		
	rior		
	mode.		
	The		
	use		
	of		
	the		
	pos-		
	te-		
	rior		
	mode		
	here		
	is		
	cru-		
	cially		
	im-		
	por-		

BL8NDEDA ANTOIBLANDIEDA
tant
as
it's
the
mode
that
is
com-
pu-
ta-
tion-
ally
avail-
able
from
Ten-
sor-
Flow
and
Keras.
Ap-

11920?? proximate normality and second order approximation also holds, see ?, ? and

B20NDEDA ANIO IBILIKŅ DEDA

?

for

fu-

ture

dis-

cus-

sion.

Specif-

i-

cally,

$$\sqrt{nI(\hat{\theta}_n)} \left(\theta_n^* - \hat{\theta}_n\right) \stackrel{D}{=} Z$$

where

 $Z \sim$

N(0, 1)

is

a

stan-

dard

Nor-

mal

12120??

vari-

able.

The

con-

di-

tional

pos-

te-

rior

sat-

is-

fies

$$\mathbb{P}\left(|\theta_n^* - \hat{\theta}_n| > \epsilon\right) \to$$

0

for

each

 $\epsilon >$

0

as

 $n \rightarrow$

 ∞ .

B22NDEDA ANIOIBLUNDEDA In the 'large p'case, a number of results are available for posterior concentra-

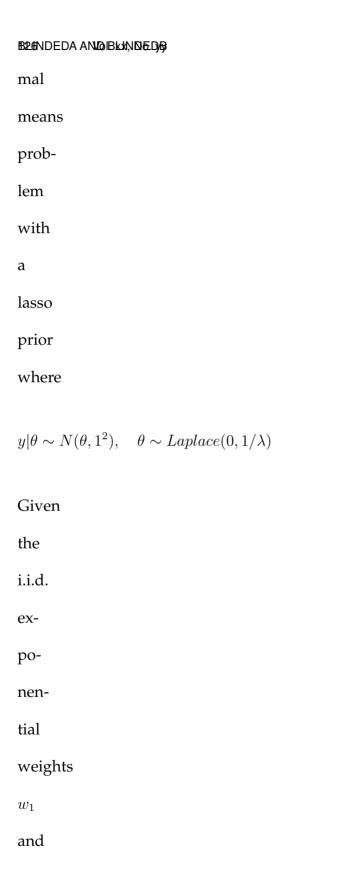
12320?? tion, for example, see ? for sparse high dimensional mod-

els.

B24NDEDA AND BLKNDEDB 3. AP-PLI-CA-TIONS Consider now a number of scenarios to assess when WBB corresponds

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

12	520??
	to
	a
	full
	Bayesian
	pos-
	te-
	rior
	dis-
	tri-
	bu-
	tion.
	3.1.
	Lasso
	First,
	a
	sim-
	ple
	uni-
	vari-
	ate
	nor-



12720?? w_2 , the weighted posterior mode $\theta^*_{\mathbf{w}}$ is $\theta_{\mathbf{w}}^* = \frac{1}{\theta \in \Theta} \left\{ \frac{w_1}{2} (y - \theta)^2 + \lambda w_2 |\theta| \right\}.$ This is sufficiently simple for

an

ex-

B28NDEDA ANIOIBkkl, NDEDA

act

WBB

so-

lu-

tion

in

terms

of

soft

thresh-

old-

ing:

$$\theta_{\mathbf{w}}^* = \begin{cases} y - \lambda w_2/w_1 & \text{if } y > \lambda w_2/w_1, \\ y + \lambda w_2/w_1 & \text{if } y < -\lambda w_2/w_1, \\ 0 & \text{if } |y| \le \lambda w_2/w_1. \end{cases}$$

The

WBB

mean

$$E_{\mathbf{w}}(\theta_{\mathbf{w}}^*|y)$$

is

12920?? approximated by the sample mean of $\{\theta_{\mathbf{w}}^{*(k)}\}_{k=1}^{K}.$ On the other hand, ? gives the expression

for

BBONDEDA ANKOIBŁKNIDEDA

the

pos-

te-

rior

mean,

$$E(\theta|y) = \frac{\int_{-\infty}^{\infty} \theta \exp\left\{-(y-\theta)^{2}/2 - \lambda|\theta|\right\} d\theta}{\int_{-\infty}^{\infty} \exp\left\{-(y-\theta)^{2}/2 - \lambda|\theta|\right\} d\theta}$$

$$= \frac{F(y)}{F(y) + F(-y)} (y+\lambda) + \frac{F(-y)}{F(y) + F(-y)} (y-\lambda)$$

$$= y + \frac{F(y) - F(-y)}{F(y) + F(-y)} \lambda$$

where

$$F(y) =$$

$$\exp(y)\Phi(-y -$$

 λ)

and

 $\Phi(\cdot)$

is

the

c.d.f.

of

13	13120??	
	stan-	
	dard	
	nor-	
	mal	
	dis-	
	tri-	
	bu-	
	tion.	
	We	
	plot	
	the	
	WBB	
	mean	
	ver-	
	sus	
	the	
	ex-	
	act	
	pos-	
	te-	
	rior	

BB2NDEDA ANTOIBKKNDEDA mean in Figure (??). Interestingly, **WBB** algorithm gives 3.2. sparser Diposbetes

Plata

Theans.

trate

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

133	320??
	our
	method-
	ol-
	ogy,
	we
	use
	weighted
	Bayesian
	Boot-
	strap
	(WBB)
	on
	the
	clas-
	sic
	di-
	a-
	betes
	dataset.
	The
	mea-
	Sure-

BB4NDEDA ANADIBŁKŅ DEDA ments for 442 diabetes patients are obtained (n =442), with 10 baseline variables (p =10),

13	13520??	
	such	
	as	
	age,	
	sex,	
	body	
	mass	
	in-	
	dex,	
	av-	
	er-	
	age	
	blood	
	pres-	
	sure,	
	and	
	six	
	blood	
	serum	
	mea-	
	sure-	
	ments.	

BB6NDEDA ANVOIBLXINDEDA



li-

hood

func-

tion

is

given

by

$$l(y|\beta) = \prod_{i=1}^{n} p(y_i|\beta)$$

where

$$p(y_i|\beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y_i - x_i'\beta)^2\right\}.$$

We

draw

1000

sets

of

weights

13	13720??			
	$\mathbf{w} =$			
	$\{w_i\}_{i=1}^{n+1}$			
	where			
	w_i 's			
	are			
	i.i.d.			
	ex-			
	po-			
	nen-			
	tials.			
	For			
	each			
	weight			
	set,			
	the			
	weighted			
	Bayesian			
	es-			
	ti-			
	mate			
	$eta_{\mathbf{w}}^*$			
	ie			

BBBNDEDA ANDIBLANDEDA calculated using (??) via the regularization method in the package glmnet.

$$\hat{\beta}_{\mathbf{w}} := \sum_{\beta=1}^{n} w_i (y_i - x_i' \beta)^2 + \lambda w_{n+1} \sum_{j=1}^{p} |\beta_j|.$$

The

reg-

u-

lar-

iza-

tion

fac-

tor

 λ

is

cho-

sen

by

cross-

validation

with

un-

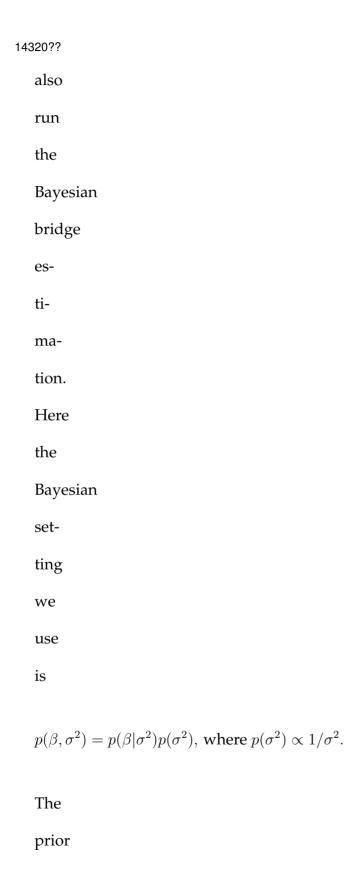
weighted

B40NDEDA ANTOIBLEN, DEDP
like-
li-
hood.
The
weighted
Bayesian
Boot-
strap
is
also
per-
formed
with
fixed
prior,
namely,
w_{n+1}
is
set
to
be

14	120??
	1
	for
	all
	boot-
	strap
	sam-
	ples.
	?
	an-
	a-
	lyze
	the
	same
	dataset
	us-
	ing
	the
	Bayesian
	Bridge
	es-
	ti-

B4L2NDEDA ANIA IBLXINIO EDA
ma-
tor
and
sug-
gest
MCMC
sam-
pling
from
the
pos-
te-
rior.
Го
com-
pare
our
WBB
re-
sults
MA

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:



B44NDEDA ANTOIBEMNISEDA on eta ,

with suit-

able

nor-

mal-

iza-

tion

con-

stant

 C_{α} ,

is

given

by

$$p(\beta) = C_{\alpha} \exp(-\sum_{j=1}^{p} |\beta_j/\tau|^{\alpha}).$$

The

hyper-

parameter

is

```
14520??
   drawn
   as
   \nu =
   	au^{-lpha} \sim
   Figuze,
   (%)here
   8hews
   the.
   re-
   sults
   of diabetes pdf
   all
   these
   three
   meth-
   ods
   (the
  weighted
   Bayesian
   Boot-
  strap
   with
```

FIGURE 2: Diabetes ex-

weighted Bayesian Boot-

strap (with

fixed prior and

weighted prior) and Bayesian Bridge

are

used to draw

from the

riors
for β_j 's,
j
=

marginal poste-

1,2,...10.

ample: the

DOI: The Canadian Journal of Statistics / La revue canadienne de statistique

B46NDEDA AND BLXNDEDB fixed prior weighted prior and the Bayesian Bridge). Marginal posteriors for β_j 's are presented. One no-

14	720??
	table
	fea-
	ture
	is
	that
	the
	weighted
	Bayesian
	Boot-
	strap
	tends
	to
	in-
	tro-
	duce
	more
	spar-
	sity
	than
	Bayesian
	Bridge

B48NDEDA ANTOIBLANDEDJB
does.
For
ex-
am-
ple,
the
weighted
Bayesian
Boot-
strap
pos-
te-
ri-
ors
of
age,
ldl
and
tch
have
higher

149	14920??		
	spikes		
	lo-		
	cated		
	around		
	0,		
	com-		
	pared		
	with		
	the		
	Bayesian		
	Bridge		
	ones.		
	For		
	tc,		
	hdl,		
	tch		
	and		
	glu,		
	multi-		
	modes		
	in		

B50NDEDA ANIO IBKKŅ DEDA the marginal posteriors are observed. In general, the posteriors with fixed priors

15120??		
	are	
	more	
	con-	
	cen-	
	trated	
	than	
	those	
	with	
	ran-	
	domly	
	weighted	
	pri-	
	ors.	
	This	
	dif-	
	fer-	
	ence	
	is	
	nat-	
	u-	
	rally	
	at-	

B52NDEDA ANIOIBLANDEDA tributed to the certainty in the prior weights. 3.3. Trend Filtering The generalized lasso solves the

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

15320??

ti-

op-

miza-

tion

prob-

lem:

$$\begin{split} \beta^* &= \int_{\beta} \left\{ l(y|\beta) + \lambda \phi(\beta) \right\} \qquad \text{(6)} \\ &= \int_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|D\beta(\beta) \| \end{split}$$

where

$$l(y|\beta) =$$

$$\frac{1}{2} || y -$$

$$X\beta\|_2^2$$

is

the

neg-

a-

tive

log-

likelihood.

 $D \in$

B54NDEDA ANDIBLANDEDA $\mathcal{R}^{m imes p}$ is a penalty matrix and $\lambda\phi(\beta) =$ $\lambda \|D\beta\|_1$ is the negative logprior or regu-

lar-

iza-

15	520??
	tion
	penalty.
	There
	are
	fast
	path
	al-
	go-
	rithms
	for
	solv-
	ing
	this
	prob-
	lem
	(see
	genlasso
	pack-
	age).

B56NDEDA AND BLXNDEDB As a subproblem, polynomial trend filtering (??) is recently introduced for piece-

15	720??	
	wise	
	poly-	
	no-	
	mial	
	curve-	
	fitting,	
	where	
	the	
	knots	
	and	
	the	
	pa-	
	ram-	
	e-	
	ters	
	are	
	cho-	
	sen	
	adap-	
	tively.	
	In-	

B58NDEDA AND BLKNDEDB tuitively, the trendfiltering estimator is similar to an adaptive spline model: it pe-

15920??		
	nal-	
	izes	
	the	
	dis-	
	crete	
	deriva-	
	tive	
	of	
	or-	
	der	
	k,	
	re-	
	sult-	
	ing	
	in	
	piece-	
	wise	
	poly-	
	no-	
	mi-	
	als	

B60NDEDA ANIO IBŁXI, NO.E.D/B of higher degree for larger k. Specifically, X = I_p in the trend filtering setting and the

16	120??
	data
	y =
	$(y_1,,y_p)$
	are
	as-
	sumed
	to
	be
	mean-
	ing-
	fully
	or-
	dered
	from
	1
	to
	p.
	The
	penalty
	ma-
	trix

B62NDEDA ANVOIBLANDEDA

is

spe-

cially

de-

signed

by

the

dis-

crete

(k +

1)-

th

or-

der

deriva-

tive,

$$D^{(1)} = \begin{bmatrix} -1 & 1 & 0 & \dots & 0 & 0 \\ 0 & -1 & 1 & \dots & 0 & 0 \\ & & & & & & \\ 0 & 0 & 0 & \dots & -1 & 1 \end{bmatrix}_{(p-1) \times p}$$

16320??
and
$D^{(k+1)} =$
$D^{(1)}D^{(k)}$
for
k =
1, 2, 3
For
ex-
am-
ple,
the
log-
prior
in
lin-
ear
trend
fil-

ter-

ing

is

B64NDEDA ANADIBŁXŅDEDA

ex-

plic-

itly

writ-

ten

as

$$\lambda \sum_{i=1}^{p-2} |\beta_{i+2} - \beta_{i+2}|$$

$$2\beta_{i+1} +$$

 β_i |.

For

a

gen-

eral

or-

der

k >

1,

$$||D^{(k+1)}\beta||_1 = \sum_{i=1}^{p-k-1} \left| \sum_{j=i}^{i+k+1} (-1)^{(j-i)} {k+1 \choose j-i} \beta_j \right|.$$

WBB

solves

the

fol-

low-

ing

gen-

er-

al-

ized

lasso

prob-

lem

in

each

draw:

$$\beta_{\mathbf{w}}^* = \frac{1}{\beta} \sum_{i=1}^p w_i (y_i - \beta_i)^2 + \lambda w_{p+1} \|D^{(k)}\beta\|_1$$
$$= \frac{1}{\beta} \frac{1}{2} \|Wy - W\beta\|_2^2 + \lambda \|D^{(k)}\beta\|_1$$
$$= W^{-1} \frac{1}{\tilde{\beta}} \frac{1}{2} \|\tilde{y}_{\mathbf{w}} - \tilde{\beta}_{\mathbf{w}}\|_2^2 + \lambda \|\tilde{D}_{\mathbf{w}}^{(k)}\tilde{\beta}_{\mathbf{w}}\|_1$$

B66NDEDA ANIOIBKKNDEDA

where

$$W = diag(\sqrt{w_i}/\sqrt{w_{p+1}}, ..., \sqrt{w_p}/\sqrt{w_{p+1}})$$

and

$$\tilde{y}_{\mathbf{w}} = Wy, \ \tilde{\beta}_{\mathbf{w}} = W\beta, \ \tilde{D}_{\mathbf{w}}^{(k)} = D^{(k)}W^{-1}.$$

To

il-

lus-

trate

our

method,

we

sim-

u-

late

data

 y_i

from

a

16720?? Fourier series regression $y_i = \sin\left(\frac{4\pi}{500}i\right) \exp\left(\frac{3}{500}i\right) + \epsilon_i$ for i =1, 2, ...500, where $\epsilon_i \sim$ $N(0, 2^2)$ are i.i.d. Gaussian noises. The

cu-

B68NDEDA ANDIBLANDEDA bic trend filtering result is given in Figure (??). For each i, the weighted Bayesian Boot-

16920?? strap gives a group of estimates $\{\beta^*_{\mathbf{w}}(i)\}_{j=1}^T$ where Tis the total number of draws. The

stan-

BZONDEDA ANZOIBŁKŅĪOEDA dard error of $\hat{\beta}_i$ is easily computed using these weighted bootstrap estimates.

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

17	20??
	3.4.
	Deep
	Learn-
	ing:
	MNIST
	Ex-
	am-
	ple
	Deep
	learn-
	ing
	is
	a
	form
	of
	ma-
	chine
	learn-
	ing
	that
	uses
	hi-

BZ2NDEDA ANZOIBŁKŅ ZIEDJĘ erarchical abstract layers of latent variables to perform pattern matching and pre-

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

17	320??
	dic-
	tion.
	?
	take
	a
	Bayesian
	prob-
	a-
	bilis-
	tic
	per-
	spec-
	tive
	and
	pro-
	vide
	a
	num-
	ber
	of
	in-

BZ4NDEDA AND BŁKŅDEDB sights into more efficient algorithms for optimization and hyperparameter tuning.

17	520??	
	The	
	gen-	
	eral	
	goal	
	is	
	to	
	finds	
	a	
	pre-	
	dic-	
	tor	
	of	
	an	
	out-	
	put	
	y	
	given	
	a	
	high	
	di-	
	men-	

BZ6NDEDA ANIOIBŁKŅ DEDA sional input x. For a classification problem, $y \in$ $\{1, 2, ..., K\}$ is a dis-

crete

vari-

able

17	720??
	and
	can
	be
	coded
	as
	a
	<i>K</i> -
	dimensional
	0-
	1
	vec-
	tor.
	The
	model
	is
	as
	fol-
	lows.
	Let
	$z^{(l)}$
	de-

BZ8NDEDA ANZOIBŁKŅ ZIEDJĄ note the lth layer, and so x = $z^{(0)}$. The final output is the response y, which can

17920??		
	be	
	nu-	
	meric	
	or	
	cat-	
	e-	
	gor-	
	i-	
	cal.	
	A	
	deep	
	pre-	
	dic-	
	tion	
	rule	
	is	
	then	

B&ONDEDA ANVOIBLANDEDA

$$z^{(1)} = f^{(1)} \Big(W^{(0)} x + b^{(0)} \Big),$$

$$z^{(2)} = f^{(2)} \Big(W^{(1)} z^{(1)} + b^{(1)} \Big),$$

. . .

$$z^{(L)} = f^{(L)} \Big(W^{(L-1)} z^{(L-1)} + b^{(L-1)} \Big),$$

$$\hat{y}(x) = z^{(L)}.$$

Here,

 $W^{(l)}$

are

weight

ma-

tri-

ces,

and

 $b^{(l)}$

are

thresh-

old

or

18	8120??			
	ac-			
	ti-			
	va-			
	tion			
	lev-			
	els.			
	$f^{(l)}$			
	is			
	the			
	ac-			
	ti-			
	va-			
	tion			
	func-			
	tion.			
	Prob-			
	a-			
	bilis-			
	ti-			
	cally,			
	the			
	out-			

B82NDEDA ANIOIBLANDEDA put yin a classification problem is generated by a probability

model

$$p(y|x, W, b) \propto \exp\{-l(y|x, W, b)\}$$

where

$$l(y|x, W, b) =$$

$$\sum_{i=1}^{n} l_i(y_i|x_i, W, b)$$

is

the

neg-

a-

tive

cross-

entropy,

$$l_i(y_i|x_i, W, b) = l_i(y_i, \hat{y}(x_i)) = \sum_{k=1}^K y_{ik} \log \hat{y}_k(x_i)$$

where

 y_{ik}

is

0

or

1

B84NDEDA ANADIBŁKŅ DEDA and K =10. Adding the negative logprior $\lambda\phi(W,b)$, the objective function (negative log-

posterior)

to

be

min-

i-

mized

by

stochas-

tic

gra-

di-

ent

de-

scent

is

$$\mathcal{L}_{\lambda}(y,\hat{y}) = \sum_{i=1}^{n} l_i(y_i,\hat{y}(x_i)) + \lambda \phi(W,b).$$

Accordingly,

with

each

draw

B86NDEDA ANIDIBLANDEDA of weights w, WBB provides the estimates $(W_{\mathbf{w}}^*,b_{\mathbf{w}}^*)$ by solving the following optimiza-

18720?? tion problem. $(W_{\mathbf{w}}^*, b_{\mathbf{w}}^*) =_{W,b} \sum_{i=1}^n w_i l_i(y_i | x_i, W, b) + \lambda w_p \phi(W, b)$ We take the classic **MNIST** example to illustrate the

ap-

pli-

B88NDEDA ANDIBLANDEDA cation of WBB in deep learning. The **MNIST** database of handwritten digits, available from Yann Le-

 $\label{thm:condition} \textit{The Canadian Journal of Statistics/La revue canadienne de statistique} \textbf{DOI:}$

18	920??
	Cun's
	web-
	site,
	has
	60,000
	train-
	ing
	ex-
	am-
	ples
	and
	10,000
	test
	ex-
	am-
	ples.
	Here
	the
	high-
	dimensional
	x

BOONDEDA ANIZO IBILIKŅIDE.D/B is a normalized and centered fixedsize (28 × 28) image and the output \hat{y} is a

19	120??
	10-
	dimensional
	vec-
	tor,
	where
	i-
	th
	со-
	or-
	di-
	nate
	cor-
	re-
	sponds
	to
	the
	prob-
	a-
	bil-
	ity
	of

B92NDEDA ANIOIBLANIOEDA that image being the ith digit. For simplicity, we build a 2layer neural

19	320??	
	net-	
	work	
	with	
	layer	
	sizes	
	128	
	and	
	64	
	re-	
	spec-	
	tively.	
	There-	
	fore,	
	the	
	di-	
	men-	
	sions	
	of	
	pa-	
	ram-	
	e-	

B94NDEDA ANADIBŁKŅDEDA

ters

are

$$W^{(0)} \in \mathcal{R}^{128 \times 784}, b^{(0)} \in \mathcal{R}^{128},$$

$$W^{(1)} \in \mathcal{R}^{64 \times 128}, \ b^{(1)} \in \mathcal{R}^{64},$$

$$W^{(2)} \in \mathcal{R}^{10 \times 64}, \ b^{(0)} \in \mathcal{R}^{10}.$$

The

ac-

ti-

va-

tion

func-

tion

 $f^{(i)}$

is

ReLU,

$$f(x) =$$

 $\max\{0,x\},$

and

the

neg-

19520?? ative logprior is specified as $\lambda \phi(W, b) = \lambda \sum_{l=0}^{2} \|W^{(l)}\|_{2}^{2}$ where $\lambda =$ 10^{-4} . Figure (??) shows the

pos-

te-

B96NDEDA ANIA IBLKNIDEDA
rior
dis-
tri-
bu-
tion
of
the
clas-
si-
fi-
са-
tion
ac-
cu-
racy
in
the
test
dataset.
We
see

1	9720??
	that
	the
	test
	ac-
	cu-
	ra-
	cies
	are
	cen-
	tered
	around
	0.75
	and
	the
	pos-
	te-
	rior
	dis-
	tri-
	bu-
	tion

B928NDEDA ANADIBŁKŅ DEDA is leftskewed. Furthermore, the accuracy is higher than 0.35 in 99% of the cases. The 95%

19920?? inter-Mag-**G**US-819N7, M&9ghted Bayesian Bootstrap (WBB) provides a computationally attractive so-

DOI: The Canadian Journal of Statistics / La revue canadienne de statistique

Figure 4:

Posterior distribu-

tion of

the classifi-

cation

accuracy.

 $n = 500, \lambda = 10^{-4}.$

BOONDEDA ANVOIBLANN DEDA
lu-
tion
to
scal-
able
Bayesian
in-
fer-
ence
(??)
whilst
ac-
count-
ing
for
pa-
ram-
e-
ter
un-
cer-

20	120??
	tainty
	by
	draw-
	ing
	sam-
	ples
	from
	a
	weighted
	pos-
	te-
	rior
	dis-
	tri-
	bu-
	tion.
	WBB
	can
	also
	be
	used

BO2NDEDA ANVOIBLAN DEDA in conjunction with proximal methods (?, ?) to provide sparsity in high dimen-

203	20??	
	sional	
	sta-	
	tis-	
	tica	
	prob-	
	lems.	
	With	
	a	
	sim-	
	i-	
	lar	
	ease	
	of	
	com-	
	pu-	
	ta-	
	tion,	
	WBB	
	pro-	
	vides	
	an	

BO4NDEDA ANVOIBLANDED/G
al-
ter-
na-
tive
to
ABC
meth-
ods
(?)
and
Vari-
a-
tional
Bayes
(VB)
meth-
ods.
A
fruit-
ful
area

205	520??
	for
	fu-
	ture
	re-
	search
	is
	the
	com-
	par-
	i-
	son
	of
	ар-
	prox-
	i-
	mate
	Bayesian
	com-
	pu-
	ta-
	tion

BOONDEDA ANDOIBLANDEDA
with
sim-
u-
lated
Bayesian
Boot-
strap
in-
fer-
ence.
ACKNOWLEDGEMENTS
Place
all
ac-
knowl-
edge-
ments
here.
In
your
ini-

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

207	20720??	
	tial	
	and	
	re-	
	vised	
	sub-	
	mis-	
	sion,	
	en-	
	sure	
	that	
	any	
	ac-	
	knowl-	
	edge-	
	ments	
	are	
	anony-	
	mous;	
	in-	
	clude	
	the	

BOSNDEDA AND IBLANDEDA full acknowledgements only after your paper has been accepted. Granting agencies should not

209	20920??				
	be				
	ab-				
	bre-				
	vi-				
	ated,				
	and				
	do				
	not				
	in-				
	clude				
	grant				
	num-				
	bers.				
	We				
	are				
	grate-				
	ful				
	for				
	your				
	as-				
	sis-				

BLONDEDA ANIABLININE.DA
tance
with
our
pub-
li-
ca-
tion
pro-
cess.
APPENDIX
AP-
PENDIXSTOCHASTIC
GRA-
DI-
ENT
DE-
SCENT
(SGD)
Stochastic
gra-

The Canadian Journal of Statistics / La revue canadienne de statistiqueDOI:

21	20??
	di-
	ent
	de-
	scent
	(SGD)
	method
	or
	its
	vari-
	a-
	tion
	is
	typ-
	i-
	cally
	used
	to
	find
	the
	deep
	learn-

BL2NDEDA ANIOIBKKNIDED/B ing model weights by minimizing the penalized loss function, $\sum_{i=1}^{n} w_i l_i(y_i; \theta) +$ $\lambda w_p \phi(\theta)$. The method mini-

213	21320??	
	mizes	
	the	
	func-	
	tion	
	by	
	tak-	
	ing	
	a	
	neg-	
	a-	
	tive	
	step	
	along	
	an	
	es-	
	ti-	
	mate	
	g^k	
	of	
	the	
	gra-	

BL4NDEDA ANADIBŁKŅ DEDA dient $\nabla \left[\sum_{i=1}^{n} w_i l_i(y_i; \theta^k) + \lambda w_p \phi(\theta^k) \right]$ at iteration k. The approximate gradient is

es-

ti-

mated

21520?? by calculating $g^{k} = \frac{n}{b_{k}} \sum_{i \in E_{k}} w_{i} \nabla l_{i}(y_{i}; \theta^{k}) + \lambda w_{p} \frac{n}{b_{k}} \nabla \phi(\theta^{k})$ Where $E_k \subset$ $\{1,\ldots,n\}$ and $b_k =$ $|E_k|$ is the number of

el-

e-

ments

BLENDEDA ANLOIBLANDEDA in E_k . When $b_k >$ 1 the algorithm is called batch **SGD** and simply SGD otherwise. A

217	21720??	
	usual	
	strat-	
	egy	
	to	
	choose	
	sub-	
	set	
	E	
	is	
	to	
	go	
	cycli-	
	cally	
	and	
	pick	
	con-	
	sec-	
	u-	
	tive	
	el-	
	٥-	

BL&NDEDA ANLOIBLAN, DEDA ments of $\{1,\ldots,T\}$, $E_{k+1} =$ $[E_k]$ $\mod n] +$ 1. The approximated direction

 g^k

is

cal-

cu-

lated

us-

219	21920??	
	ing	
	a	
	chain	
	rule	
	(aka	
	back-	
	propagation)	
	for	
	deep	
	learn-	
	ing.	
	It	
	is	
	an	
	un-	
	bi-	
	ased	
	es-	
	ti-	
	ma-	
	tor.	

B20NDEDA ANIO IBKKŅ DEDA Thus, at each iteration, the **SGD** updates the solution $\theta^{k+1} = \theta^k - t_k g^k$

ing

For

deep

learn-

22120?? applications the step size t_k (a.k.a learning rate) is usually kept constant or some simple

B22NDEDA ANIOIBKKNDED/B step size reduction strategy is used, $t_k =$ $a\exp(-kt)$. Арpropriate learning rates or the hy-

223	22320??	
	per-	
	pa-	
	ram-	
	e-	
	ters	
	of	
	re-	
	duc-	
	tion	
	sched-	
	ule	
	are	
	usu-	
	ally	
	found	
	em-	
	pir-	
	i-	
	cally	
	from	
	nu-	
	mer-	

B224INDEDA ANIOIBIAN, DEDA ical experiments and observations of the loss function progression.

22520??		
Received 9		
Au-		
gust		
2018		
Accepted 8		
Septem-		
ber		
2018		