

Online Learning and Bregman Divergences

Gunnar Rätsch[&] and Manfred Warmuth^{\$}

[&] *Australian National University, Canberra, Australia*

^{\$} *University of California at Santa Cruz, USA*

[&] <http://mlg.anu.edu.au/~raetsch>

^{\$} <http://www.cse.ucsc.edu/~manfred>

Help with the tutorial: Claudio Gentile, Jyrki Kivinen

Content of this tutorial

- L 1: Introduction to Online Learning
 - The Learning setting
 - Predicting as good as the best expert
 - Predicting as good as the best linear combination of experts
- L 2: Bregman divergences and Loss bounds
 - Introduction to Bregman divergences
 - Relative loss bounds for the linear case
 - Nonlinear case & matching losses
 - Duality and relation to exponential families
- L 3: Extensions, interpretations, applications
 - Online to Batch Conversions
 - Prior information on the weight vector
 - Some applications

Goal: How to learn efficiently in an online setting? How can one prove it?

Sources of Information

Internet <http://mlg.anu.edu.au/~raetsch/online>

Journals and papers Machine Learning, Journal of Machine Learning Research, Information and Computation, IEEE Transaction on Information Theory, ...

Small collection at <http://mlg.anu.edu.au/~raetsch/online>

Conferences Neural Information Processing Systems (NIPS), Computational Learning Theory (COLT)

People (incomplete) List at <http://mlg.anu.edu.au/~raetsch/online>

Software Not necessary, algorithms are too simple

The Online Learning setting

Protocol:

For $t = 1$ To T Do

Get instance	$\mathbf{x}_t \in \mathbf{R}^n$
Predict	$\hat{y}_t \in \mathcal{Y}$
Get label	$y_t \in \mathcal{Y}$
Incur loss	$L(y_t, \hat{y}_t, \mathbf{x}_t)$

Problem instances:

- Classification: $\hat{y}_t, y_t \in \{0, 1\}$, e.g. $L(y, \hat{y}, \mathbf{x}) = |y - \hat{y}|$
- Regression: $\hat{y}_t, y_t \in \mathbf{R}$, e.g. $L(y, \hat{y}, \mathbf{x}) = (y - \hat{y})^2$
- Density estimation: no label, e.g. $L(y, \hat{y}, \mathbf{x}) = -\log P(\mathbf{x}|\theta)$

Goal: small total loss $\sum_t L(y_t, \hat{y}_t, \mathbf{x}_t)$

Predicting almost as good as the best expert

	E_1	E_2	E_3	\dots	E_n	prediction	<i>true label</i>	loss
day 1	1	1	0	\dots	0	0	1	1
day 2	1	0	1	\dots	0	1	0	1
day 3	0	1	1	\dots	1	1	1	0
day t	$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	\dots	$x_{t,n}$	\hat{y}_t	y_t	$ y_t - \hat{y}_t $

Master Algorithm

For $t = 1$ To T Do

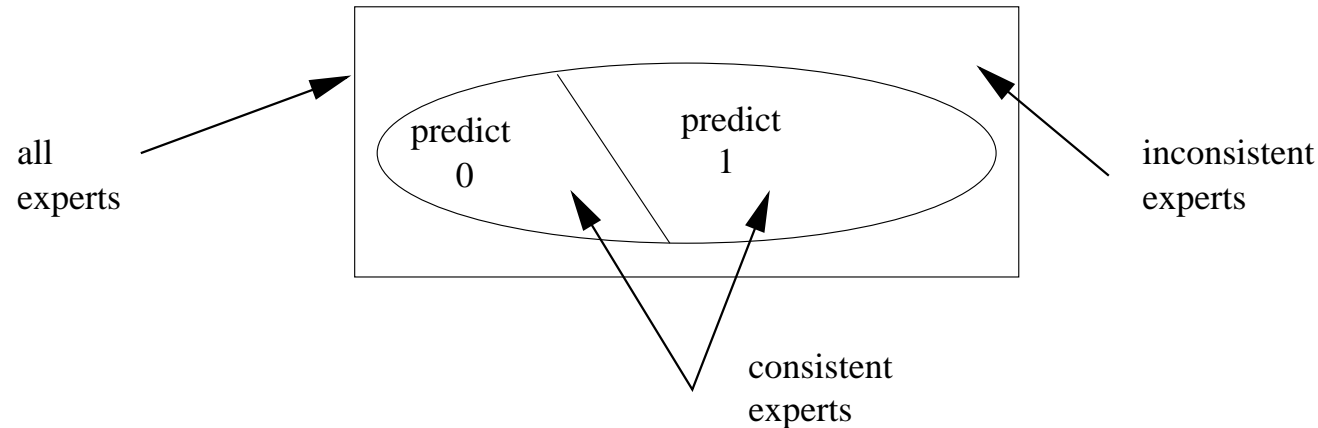
Get instance $\mathbf{x}_t \in \{0, 1\}^n$

Predict $\hat{y}_t \in \{0, 1\}$

Get label $y_t \in \{0, 1\}$

Incur loss $|y_t - \hat{y}_t|$

Halving Algorithm [BF]



- Predicts with majority
- If mistake is made, then number of **consistent** experts is (at least) **halved**
- any mistake is “converted” into knowledge on the learning problem: *mistake driven learning*

A run of the Halving Algorithm

E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	<i>majority</i>	<i>true label</i>	loss
1	1	0	0	1	1	0	0	1	0	1
x	x	0	1	x	x	1	1	1	1	0
x	x	x	1	x	x	0	0	0	1	1
x	x	x	↑	x	x	x	x			
			<i>consistent</i>							

For any sequence with a *consistent* expert

Halving Algorithm makes at most $\leq \log_2 n$ mistakes

(Good bound, but not optimal)

What if no expert is consistent?

Sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

- **total loss** of algorithm A : $L_A(S) = \sum_{t=1}^T L(A(\mathbf{x}_t), y_t)$
- **total loss** of i -th expert E_i : $L_i(S) = \sum_{t=1}^T L(E_i(\mathbf{x}_t), y_t)$

Want bounds of the form:

$$\forall S : L_A(S) \leq a \min_i L_i(S) + b \log(n)$$

where a, b are constants

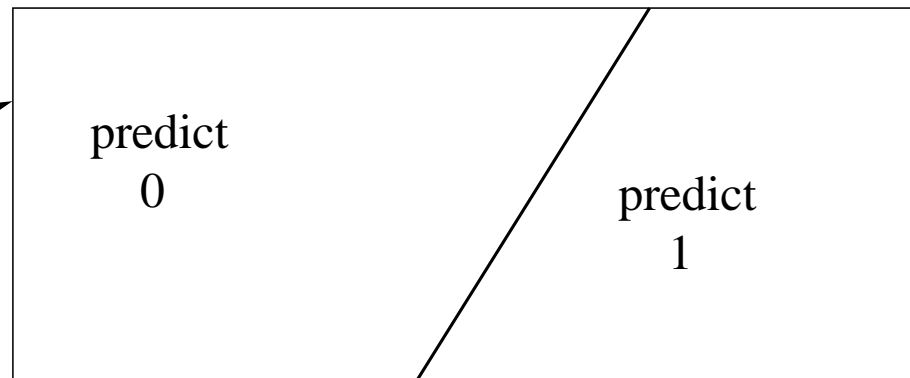
Bounds loss of algorithm **relative to** loss of best expert

Weighted Majority Algorithm [LW]

Can't wipe out experts!

Need one weight per expert

all
experts
vote with
their weight



- Predicts with larger side
- Weights of wrong experts are **multiplied** by $\beta \in [0, 1)$

Number of mistakes of the WM algorithm

$M_{t,i}$ = # of mistakes of E_i before trial t

$w_{t,i}$ = $\beta^{M_{t,i}}$ weight of E_i at beginning of trial t

W_t = $\sum_{i=1}^n w_{t,i}$ total weight at trial t

Minority $\leq \frac{1}{2}W_t$, Majority $\geq \frac{1}{2}W_t$

If no mistake then minority multiplied by β :

$$W_{t+1} \leq 1 W_t$$

If mistake then majority multiplied by β :

$$W_{t+1} \leq \underbrace{1}_{\text{minority}} \frac{1}{2}W_t + \underbrace{\beta}_{\text{majority}} \frac{1}{2}W_t = \frac{1+\beta}{2} W_t$$

M = Total number of mistakes of WM

M_i = Total number of mistakes of expert E_i

Hence

$$\underbrace{W_{T+1}}_{\text{total final weight}} \leq \left(\frac{1 + \beta}{2} \right)^M W_1$$

$$W_{T+1} = \sum_{j=1}^n w_{T+1,j} = \sum_{j=1}^n \beta^{M_{T+1,j}} \geq \beta^{M_{T+1,i}}$$

where $i = \operatorname{argmin}_j M_{T+1,j}$

We got:

$$\left(\frac{1 + \beta}{2} \right)^M \underbrace{W_1}_n \geq \beta^{M_i}$$

Relative Loss bound for Weighted Majority

Solving for M :

$$M \leq \frac{\ln \frac{1}{\beta}}{\ln \frac{2}{1+\beta}} M_i + \frac{1}{\ln \frac{2}{1+\beta}} \ln n$$

$$M \leq \underbrace{2.63}_a \min_i M_i + \underbrace{2.63}_b \ln n$$

$\beta = 1/e$

For **all** sequences, loss of master algorithm
is **comparable to** loss of best expert
 \Rightarrow **Relative** loss bounds

[Fr]

Other Loss Functions

absolute loss $L(y, \hat{y}) = |y - \hat{y}|$

square loss $L(y, \hat{y}) = (y - \hat{y})^2$

entropic loss $L(y, \hat{y}) = y \ln \frac{y}{\hat{y}} + (1 - y) \ln \frac{1-y}{1-\hat{y}},$
 $y, \hat{y} \in [0, 1]$

entropic loss \pm $L(y, \hat{y}) = \frac{1+y}{2} \ln \frac{1+y}{1+\hat{y}} + \frac{1-y}{2} \ln \frac{1-y}{1-\hat{y}},$
 $y, \hat{y} \in [-1, +1]$

hellinger loss $L(y, \hat{y}) = \frac{1}{2} (\sqrt{1-y} - \sqrt{1-\hat{y}})^2 + \frac{1}{2} (\sqrt{y} - \sqrt{\hat{y}})^2$
 $y, \hat{y} \in [0, 1]$

How does it work with other loss functions?

One weight per expert:

[V]

$$w_{t,i} = \beta^{L_{t,i}} = e^{-\eta L_{t,i}}$$

where $L_{t,i}$ is total loss of E_i before trial t
and η is a positive learning rate

Master predicts with the **weighted average** (WA)

[KW]

$$v_{t,i} = \frac{w_{t,i}}{\sum_{i=1}^n w_{t,i}} \quad \text{normalized weights}$$
$$\hat{y}_t = \sum_{i=1}^n v_{t,i} x_{t,i} = \mathbf{v}_t \cdot \mathbf{x}_t$$

where $x_{t,i}$ is the prediction of E_i in trial t

Bounds for other Loss Functions

\forall sequences S of examples $\langle (\mathbf{x}_t, y_t) \rangle_{1 \leq t \leq T}$
 where $\mathbf{x}_t \in [0, 1]^n$ and $y_t \in [0, 1]$

$$L_{WA}(S) \leq \min_i \underbrace{1}_a L_i(S) + \underbrace{1/\eta}_b \ln(n)$$

$1/\eta$	WA	fancy
entropic	1	1
square	2	1/2
hellinger	1	.71

- Improved constants of $1/\eta$ when Master uses fancier pred. [V]
- For the discrete loss and the absolute loss: $a > 1$

Proof

$$\begin{aligned} \text{Key inequality: } L(y, \mathbf{v}_t \cdot \mathbf{x}_t) &\leq \frac{1}{\eta} \ln W_t - \frac{1}{\eta} \ln W_{t+1} \\ &= -\frac{1}{\eta} \ln \frac{W_{t+1}}{W_t} \end{aligned}$$

Telescoping:

$$\begin{aligned} L_{\text{WA}}(S) &\leq -\frac{1}{\eta} \ln \frac{W_{T+1}}{W_1} \\ &= -\frac{1}{\eta} \ln \sum_{j=1}^n \frac{1}{n} e^{-\eta L_{T+1,j}(S)} \\ &\leq -\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_i(S)} && i = \operatorname{argmin}_j L_{T+1,j} \\ &= -\frac{1}{\eta} \ln \frac{1}{n} e^{-\eta L_{T+1,i}(S)} \\ &= L_{T+1,i}(S) + \frac{1}{\eta} \ln n \end{aligned}$$

Usefulness:

- Easy to combine many pretty good experts (algorithms) so that Master is guaranteed to be almost as good as the best
- Bounds **logarithmic** in number of experts (**multiplicative** updates)

Questions:

- How to obtain algorithms that do well compared to best **linear combination** or best **thresholded linear combination** of experts?
- How to motivate the updates?
- What are good measures of progress?
- What are good loss functions?
- Methods for proving relative loss bounds?

A more general setting

Example	Prediction of alg A	Label	Loss of alg A
x_1	\hat{y}_1	y_1	$L(y_1, \hat{y}_1)$
\vdots	\vdots	\vdots	\vdots
x_t	\hat{y}_t	y_t	$L(y_t, \hat{y}_t)$
\vdots	\vdots	\vdots	\vdots
x_T	\hat{y}_T	y_T	$L(y_T, \hat{y}_T)$
		Total Loss	$L_A(S)$

Sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T)$

Comparison class $\{\mathbf{u}\}$

Relative loss $L_A(S) - \inf_{\{\mathbf{u}\}} L_{\mathbf{u}}(S)$

Goal: Bound relative loss for arbitrary sequence of examples

Example: Learning Disjunctions of Experts

variables/experts						
E_1	E_2	E_3	E_4	<i>true label</i>	$E_1 \vee E_3$	$E_3 \vee E_4$
1	1	0	0	0	1	0
1	0	1	0	1	1	1
0	1	1	1	0	1	1
0	1	0	0	1	0	0
$x_{t,1}$	$x_{t,2}$	$x_{t,3}$	$x_{t,4}$		↑	↑
					3	2
					mistakes	

$E_1 \vee E_3$ becomes $\mathbf{u} = (1, 0, 1, 0)$

$E_1 \vee E_3$ is one on $\mathbf{x}_t \in \{0, 1\}^n$ iff $\mathbf{u} \cdot \mathbf{x}_t \geq 1$

Weighted Majority on k -literal Disjunctions

Do as well as best k out of n literal (monotone) disjunction

Each disjunction is an expert

Keep one weight per disjunction: $\binom{n}{k}$ weights

$$\begin{array}{l} \# \text{ of mistakes} \\ \text{of WM} \end{array} \leq 2.63 M + 2.63 k \ln \frac{n}{k}$$

M is # of mistakes of best

Time (and space) **exponential** in k

Efficient algorithm have only one weight per **literal**

The Perceptron Algorithm

In trial t : Get instance $\mathbf{x}_t \in \{0, 1\}^n$

If $\mathbf{w}_t \cdot \mathbf{x}_t \geq 1/2$ then $\hat{y}_t = 1$

else $\hat{y}_t = 0$

Get label $y_t \in \{0, 1\}$

If **mistake** then

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta (\hat{y}_t - y_t) \mathbf{x}_t$$

k -literal Disjunctions with Perceptron

Perceptron Convergence Theorem ($\eta = \frac{1}{2n}$)

$$\# \text{ of mistakes} \leq 4A + 4kn$$

where A is $\#$ of attribute errors of best disjunction of size k , i.e., the minimum $\#$ of attributes that need to be flipped to make the disjunction consistent

$$A \leq kM$$

Lower bound for rotation invariant algorithms:

[KWA]

$$\# \text{ mistakes} = \Omega(n)$$

The Winnow Algorithm [L]

In trial t : Get instance $\mathbf{x}_t \in \{0, 1\}^n$

If $\mathbf{w}_t \cdot \mathbf{x}_t \geq \theta$ then $\hat{y}_t = 1$

else $\hat{y}_t = 0$

Get label $y_t \in \{0, 1\}$

If **mistake** then

$$w_{t+1,i} = w_{t,i} e^{-\eta (\hat{y}_t - y_t) x_{t,i}}$$

Mistake bound ($e^{-\eta} = 1/3$, $\theta = \frac{3 \ln 3}{8}$)

[AW]

$$\# \text{ of mistakes} \leq 4A + 3.6 k \ln \frac{n}{k}$$

Not rotation invariant!

k -term DNF via Feature Expansion [KW]

$$\begin{array}{ccc}
 \mathbf{x} = & \begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} & \Rightarrow \quad \Phi(\mathbf{x}) = \begin{array}{c} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_1 x_2 \\ \vdots \\ x_1 x_2 \dots x_n \end{array} \\
 & n \text{ inputs} & 2^n \text{ features}
 \end{array}$$

k -term DNF in input space is k -literal disjunction in feature space

$$\underbrace{\Phi(\mathbf{x})}_{2^n} \cdot \underbrace{\Phi(\mathbf{y})}_{2^n} = \underbrace{\prod_{i=1}^n (1 + x_i y_i)}_{O(n) \text{ time}} = K(\underbrace{\mathbf{x}}_n, \underbrace{\mathbf{y}}_n)$$

(Simple ANOVA kernel)

k -term DNF via Feature Expansion: Perceptron

Perceptron: $\mathbf{w}_t = \sum_{q \text{ mistake}} \alpha_q \Phi(\mathbf{x}_i)$

Prediction:

$$\mathbf{w}_t \cdot \Phi(\mathbf{x}) = \left(\sum_{q \text{ mistake}} \alpha_q \Phi(\mathbf{x}_q) \right) \cdot \Phi(\mathbf{x})$$

$$= \sum_{q \text{ mistake}} \alpha_q \Phi(\mathbf{x}_q) \cdot \Phi(\mathbf{x})$$

$$= \underbrace{\sum_{q \text{ mistake}} \alpha_q K(\mathbf{x}_q, \mathbf{x})}_{\text{time: } O(n \cdot \# \text{ mistakes})}$$

Mistake bound: $O(k 2^n)$

k -term DNF via Feature Expansion: Winnow

Winnow: $w_{t,i} = \exp\left(-\eta \sum_{q \text{ mistake}} \alpha_q \Phi(\mathbf{x}_q)_i\right)$

log of weights is linear combination of past examples

Mistake bound: $O(k \ln 2^n) = O(kn)$

prediction time: $\Omega(2^n \# \text{ mistakes})$

No kernel trick with purely multiplicative updates!

So far

- Learning relative to best expert and best disjunction
- Various loss functions
- Perceptron versus Winnow and expansion into feature space

On-line Linear Regression

For $t = 1, \dots, T$ do

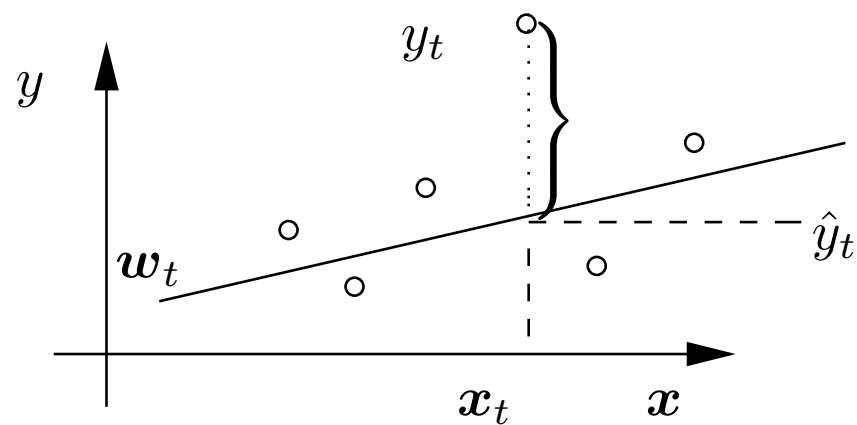
Get instance $\mathbf{x}_t \in \mathbf{R}^n$

Predict $\hat{y}_t = \mathbf{w}_t \cdot \mathbf{x}_t$

Get label $y_t \in \mathbf{R}$

Incur loss $L_t(\mathbf{w}_t) = (y_t - \hat{y}_t)^2$

Update \mathbf{w}_t to \mathbf{w}_{t+1}



Assume comparison class $\{u\}$ is a set of **linear** predictors

$$u : x \rightarrow u \cdot x$$

Examples of Updates

Gradient descent

($\mathbf{w} \in \mathbf{R}^n$)

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \nabla L_t(\mathbf{w}_t) \\ &= \mathbf{w}_t - \eta (\mathbf{w}_t \cdot \mathbf{x}_t - y_t) \mathbf{x}_t\end{aligned}$$

[WH]

Exponentiated Gradient Algorithm

(\mathbf{w} is probability vector)

$$w_{t+1,i} = w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right] / \text{normalization}$$

[KW]

More examples of Updates

Unnormalized Exponentiated Gradient Algorithm

[KW]

($\mathbf{w} \geq \mathbf{0}$)

$$\mathbf{w}_{t+1,i} = w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]$$

Binary Exponentiated Gradient Algorithm

[By]

($\mathbf{w} \in [0, 1]^n$)

$$\mathbf{w}_{t+1,i} = \frac{w_{t,i} \exp \left[-\eta \frac{\partial N_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]}{1 - w_{t,i} + w_{t,i} \exp \left[-\eta \frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}} \right]}$$

More examples of Updates (cont)

p -norm Algorithms

[GLS, GL]

($\mathbf{w} \in \mathbf{R}^n$)

$$\mathbf{w}_{t+1} = f^{-1}(f(\mathbf{w}_t) - \eta \nabla L_t(\mathbf{w}_t))$$

where

$$f(\mathbf{w}) = \nabla \frac{1}{2} \|\mathbf{w}\|_q^2 = \nabla \frac{1}{2} \left(\sum_i |w_i|^q \right)^{2/q}$$

and q dual to p (i.e., $\frac{1}{p} + \frac{1}{q} = 1$)

- $p = 2$ becomes gradient descent
- $p = O(\log n)$ becomes EG-like algs
- $2 < p < O(\log n)$ interpolates between the two extremes

Motivation of Updates [KW]

Gradient descent

$$\begin{aligned} \mathbf{w}_{t+1} &= \underset{\mathbf{w}}{\operatorname{argmin}} \left(\|\mathbf{w} - \mathbf{w}_t\|_2^2 / 2 + \eta (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2 / 2 \right) \\ &= \mathbf{w}_t - \eta \underbrace{(\mathbf{w}_{t+1} \cdot \mathbf{x}_t - y_t)}_{\approx \mathbf{w}_t \cdot \mathbf{x}_t} \mathbf{x}_t \end{aligned}$$

Exponentiated Gradient Algorithm

$$\begin{aligned} \mathbf{w}_{t+1} &= \underset{\mathbf{w}}{\operatorname{argmin}} \left(\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}} + \eta (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2 / 2 \right) \\ &= w_{t,i} \exp \left[-\eta \underbrace{(\mathbf{w}_{t+1} \cdot \mathbf{x}_t - y_t)}_{\approx \mathbf{w}_t \cdot \mathbf{x}_t} x_{t,i} \right] / \text{normalization} \end{aligned}$$

Alternate Motivation: Continuous Updates [WJ]

Continuous time \Rightarrow Ordinary differential equations

Gradient Descent

$$\mathbf{w} \in \mathbf{R}^n$$

$$\dot{\mathbf{w}}_t = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$$

Unnormalized Exponentiated Gradient Alg.

$$\mathbf{w} \geq \mathbf{0}$$

$$\dot{\log}(\mathbf{w}_t) = -\eta \nabla_{\mathbf{w}} L_t(\mathbf{w}_t)$$

Alternate Motivation: Continuous Updates (cont)

Discretization

$$\frac{\mathbf{f}(\mathbf{w}_{t+h}) - \mathbf{f}(\mathbf{w}_t)}{h} = -\eta \nabla L_t(\mathbf{w}_t)$$

$$\mathbf{w}_{t+h} = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_t) - \eta h \nabla \mathbf{w} L_t(\mathbf{w}_t))$$

We use $h = 1$

$$\mathbf{w}_{t+1} = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_t) - \eta_t \nabla L_t(\mathbf{w}_t))$$

Conjecture: **Forward Euler** better:

Replace $\nabla \mathbf{w} L_t(\mathbf{w}_t)$ by $\nabla \mathbf{w} L_t(\mathbf{w}_{t+h})$

Families of update algorithms

parameter “divergence”	name of family	update algorithms
$\ \mathbf{w} - \mathbf{w}_t\ _2^2$	Gradient Descent	Widrow Hoff (LMS) Linear Least Squares. Backpropagation Perceptron Algorithms kernel based algorithms,...
$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$	Exponentiated Gradient Algorithm	expert algs Normalized Winnow “AdaBoost”

Families of update algorithms (cont)

parameter “divergence”	name of family	update algorithms
$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$ $+ w_{t,i} - w_i$	Unnormalized Exp. Grad. Alg.	Winnow
$\sum_{i=1}^n w_i \ln \frac{w_i}{w_{t,i}}$ $+(1 - w_i) \ln \frac{1 - w_i}{1 - w_{t,i}}$	Binary Exp. Grad. Alg.	-
any	-	-

“Bregman divergence”

Members of different families exhibit different behavior

Additive Gradient versus Exponentiated Gradient

An experimental study: Part I

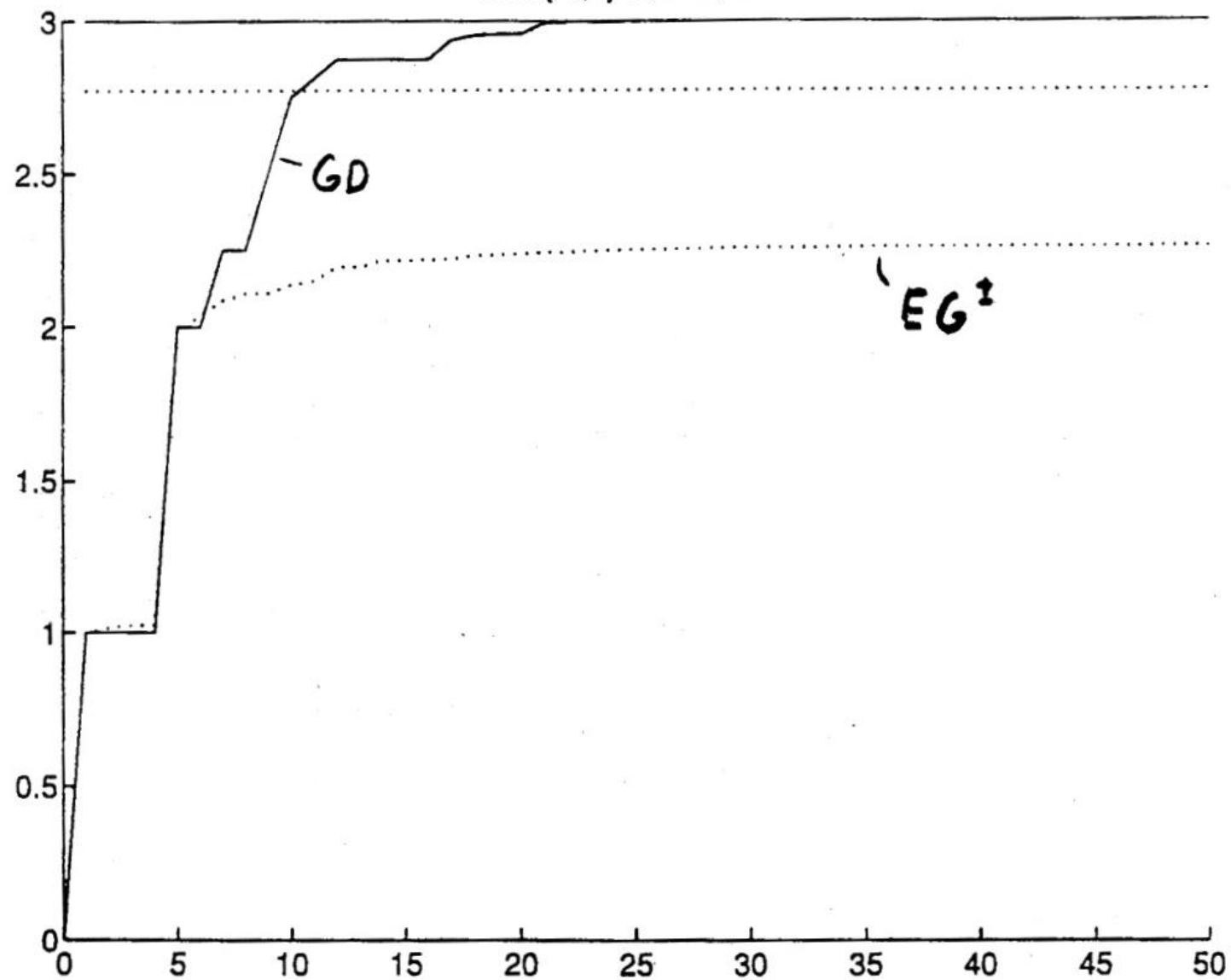
The instances are random vectors $\mathbf{x} \in \{-1, +1\}^n$

The target vector \mathbf{u} is sparse, e.g. $\mathbf{u} = [\underbrace{1, 1}_{k \text{ times}}, 0, \dots, 0]$

The labels are given $y_t = \mathbf{u} \cdot \mathbf{x}_t$

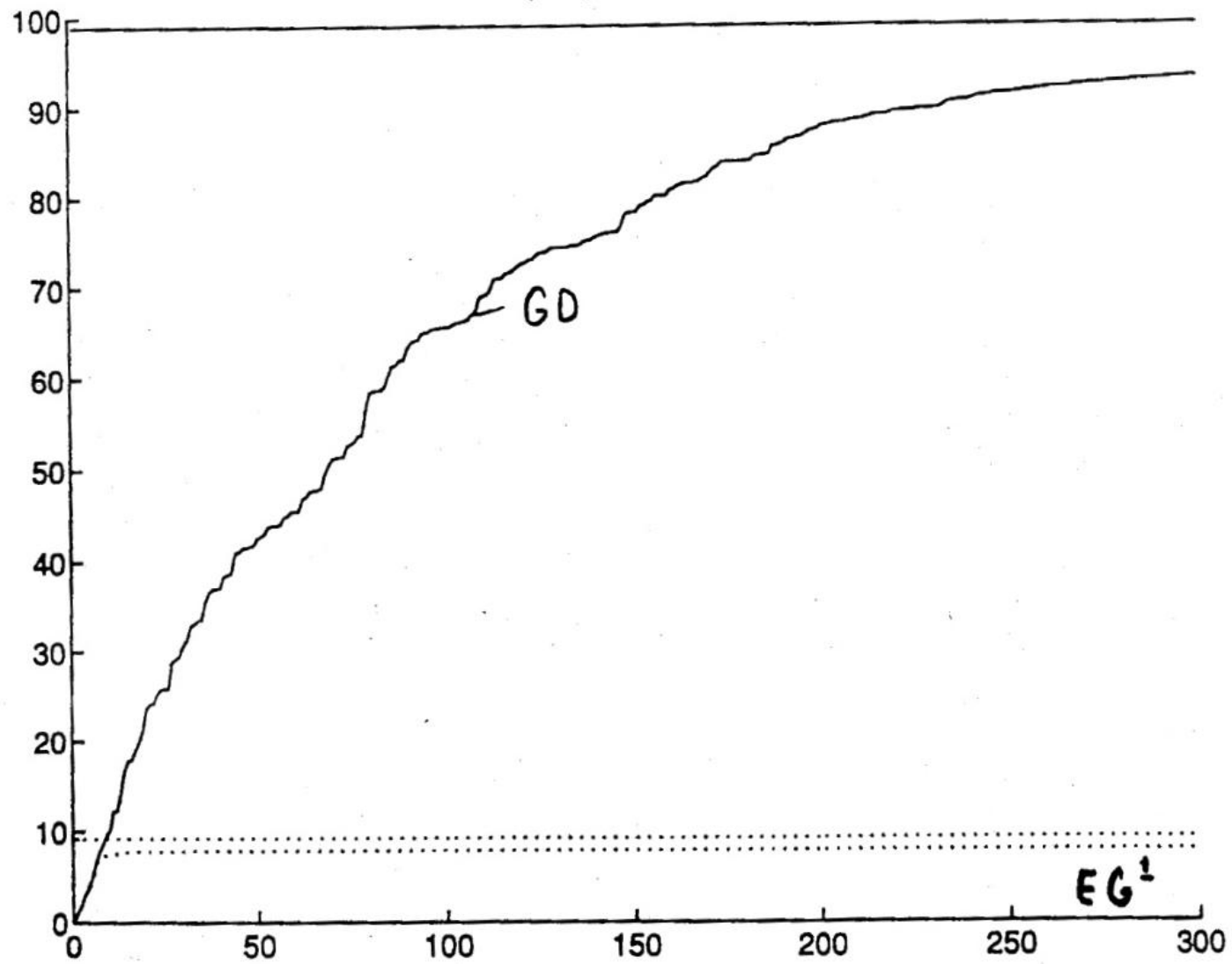
Loss function is the quadratic loss $L(y, \hat{y}) = (y - \hat{y})^2$

$x \in \{-1,1\}^n; k=1; n=4$

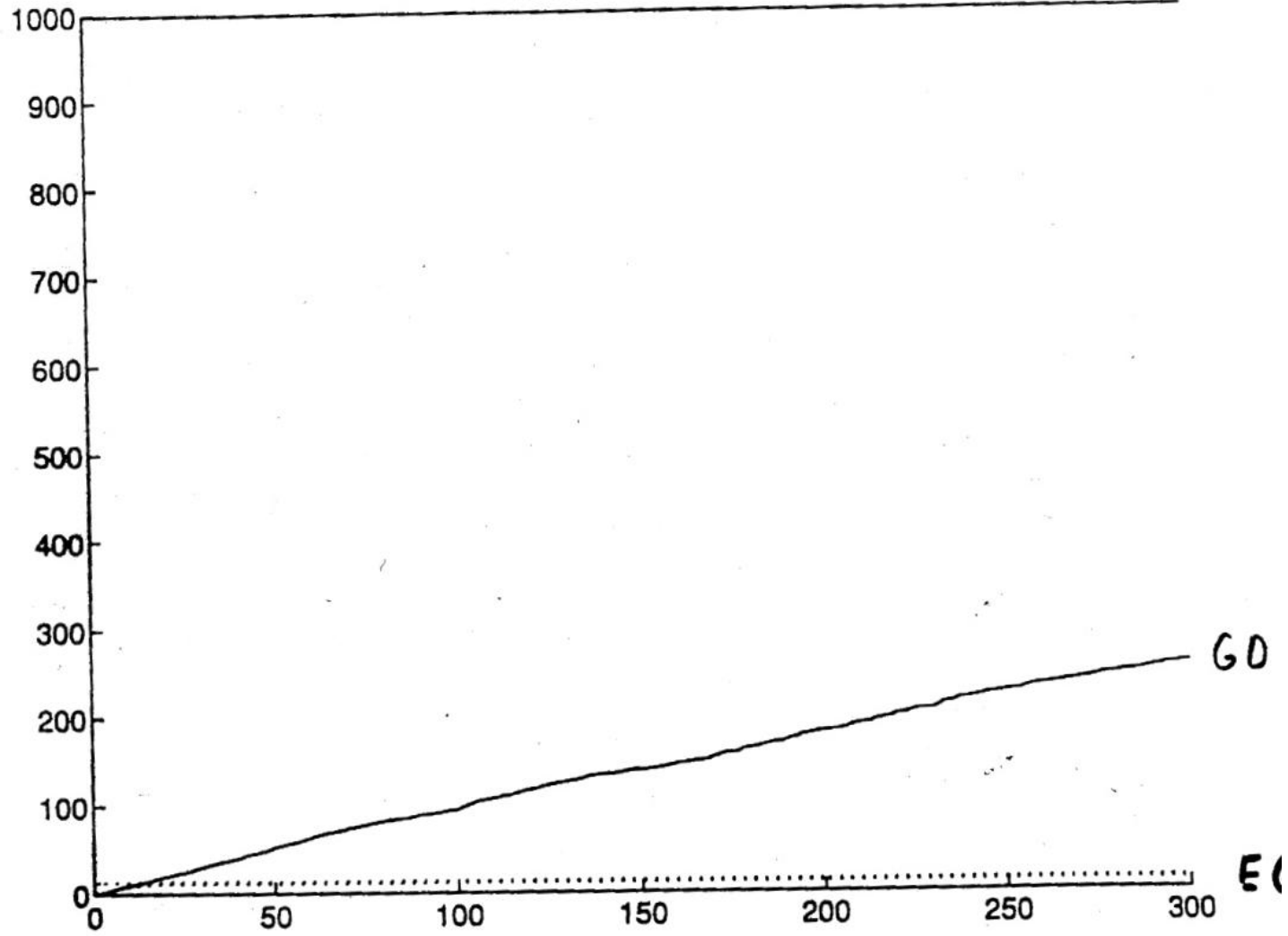


EG¹

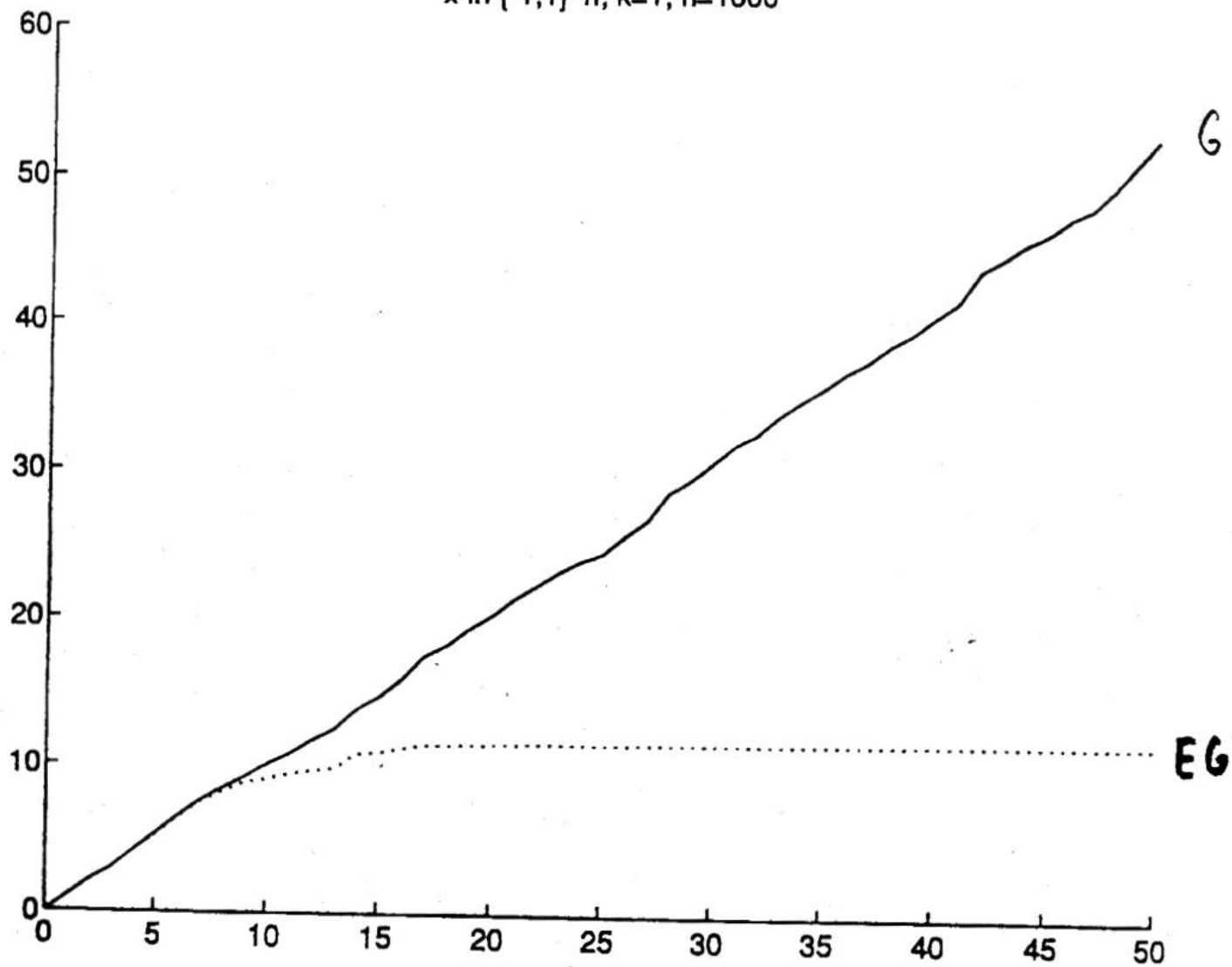
$x \in \{-1, 1\}^n; k=1; n=100$



$x \in \{-1,1\}^n; k=1; n=1000$

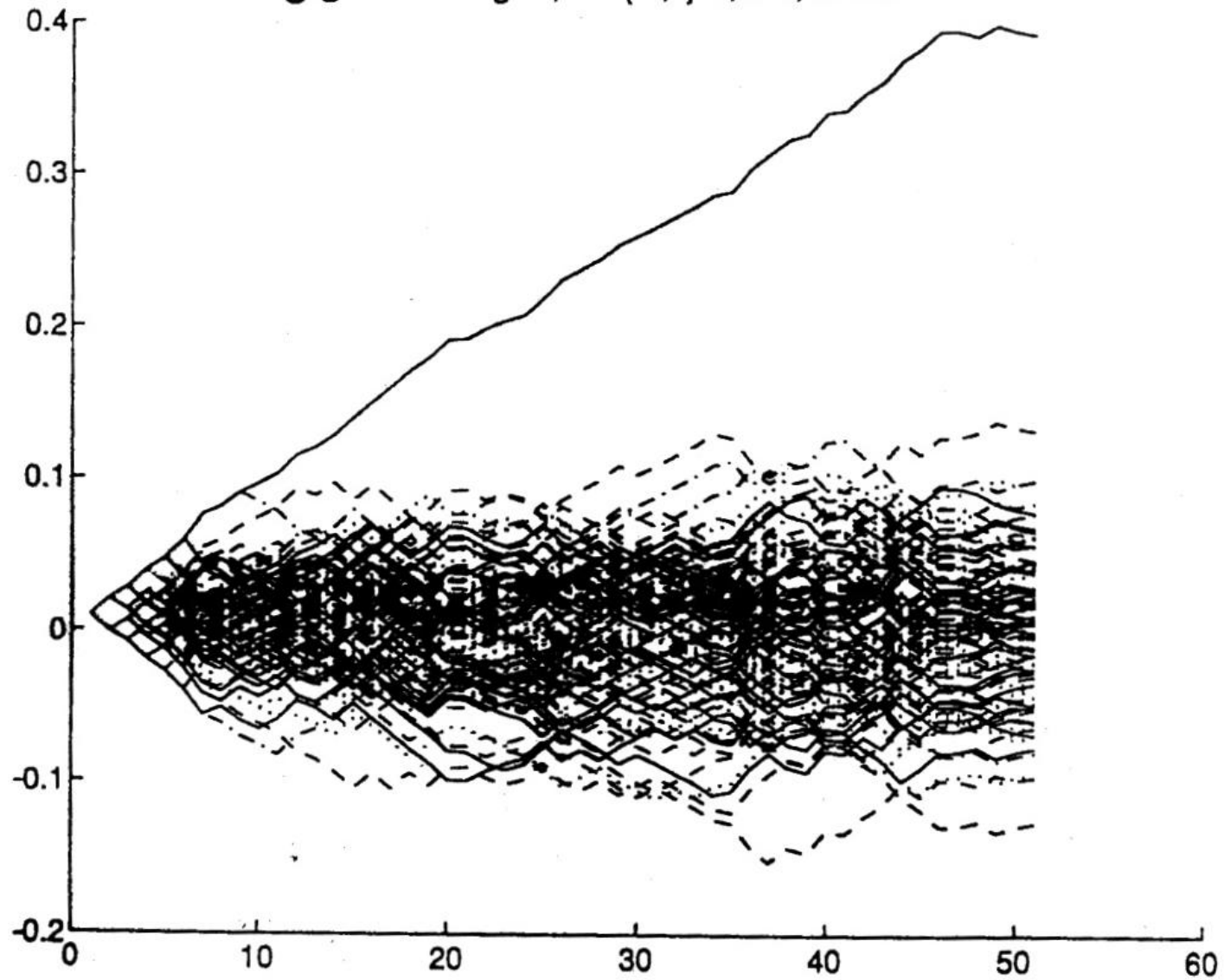


$x \text{ in } \{-1,1\}^n; k=1; n=1000$



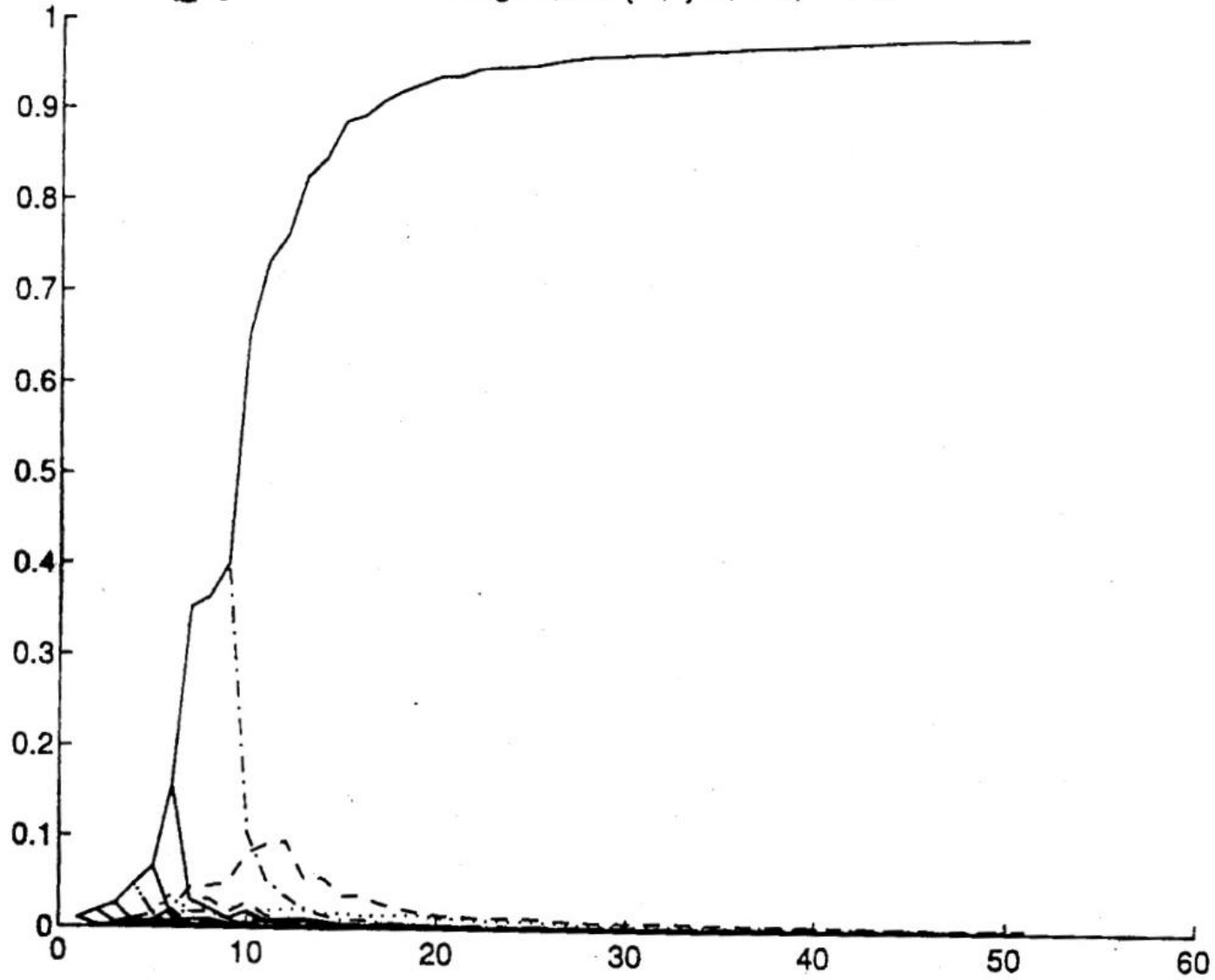
GD

weights; x in $(-1,1)^n$; $k=1$; $n=100$



EG^t

weights; x in $\{-1,1\}^n$; $k=1$; $n=100$



Loss bounds

Assume Example Sequence is

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_t, y_t), \dots \text{ where } y_t = \mathbf{u} \cdot \mathbf{x}_t$$

(the zero-error case)

Gradient Descent:

$$L_{GD}(S) \leq \left(\|\mathbf{u}\|_2 \max_t \|\mathbf{x}_t\|_2 \right)^2$$

Exponentiated Gradients:

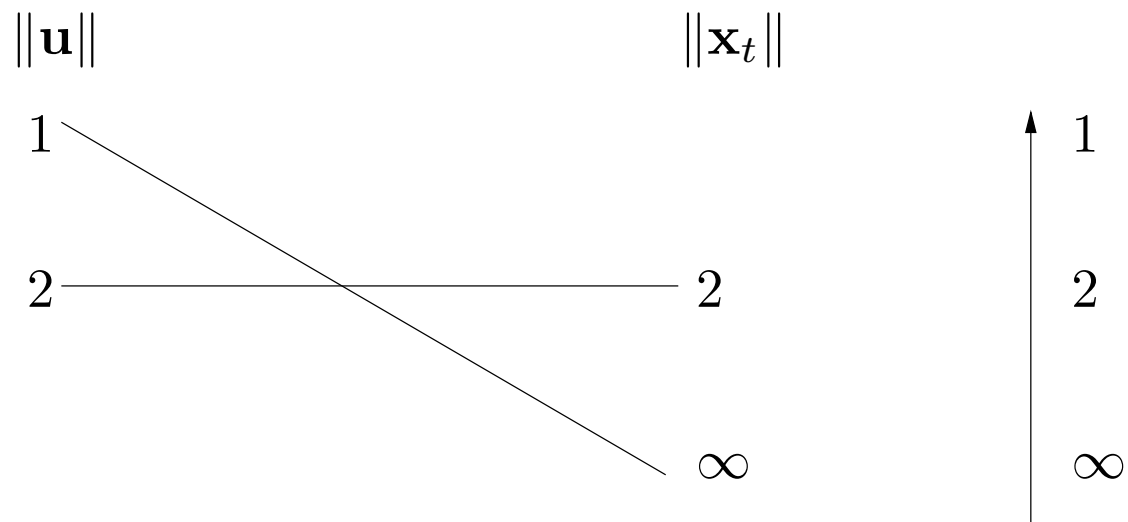
$$L_{EG\pm}(S) \leq \left(\|\mathbf{u}\|_1 \max_t \|\mathbf{x}_t\|_\infty \right)^2 \log(2n)$$

Incomparable Loss bounds

$$L_{GD}(S) \leq \left(\|\mathbf{u}\|_2 \max_t \|\mathbf{x}_t\|_2 \right)^2$$

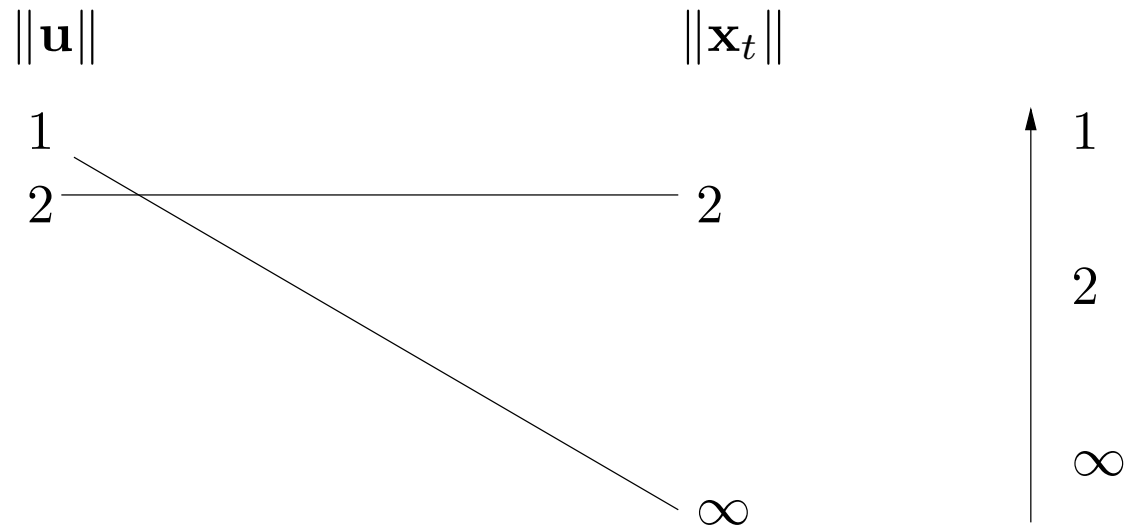
$$L_{EG\pm}(S) \leq \left(\|\mathbf{u}\|_1 \max_t \|\mathbf{x}_t\|_\infty \right)^2 \log(2n)$$

Products of two norms:



Incomparable Loss bounds (cont)

Experiments so far:



Widest gap between ℓ_2 and ℓ_∞ norms if $\mathbf{x} \in \{-1, +1\}^N$

The Worst Case: Hadamard Matrices

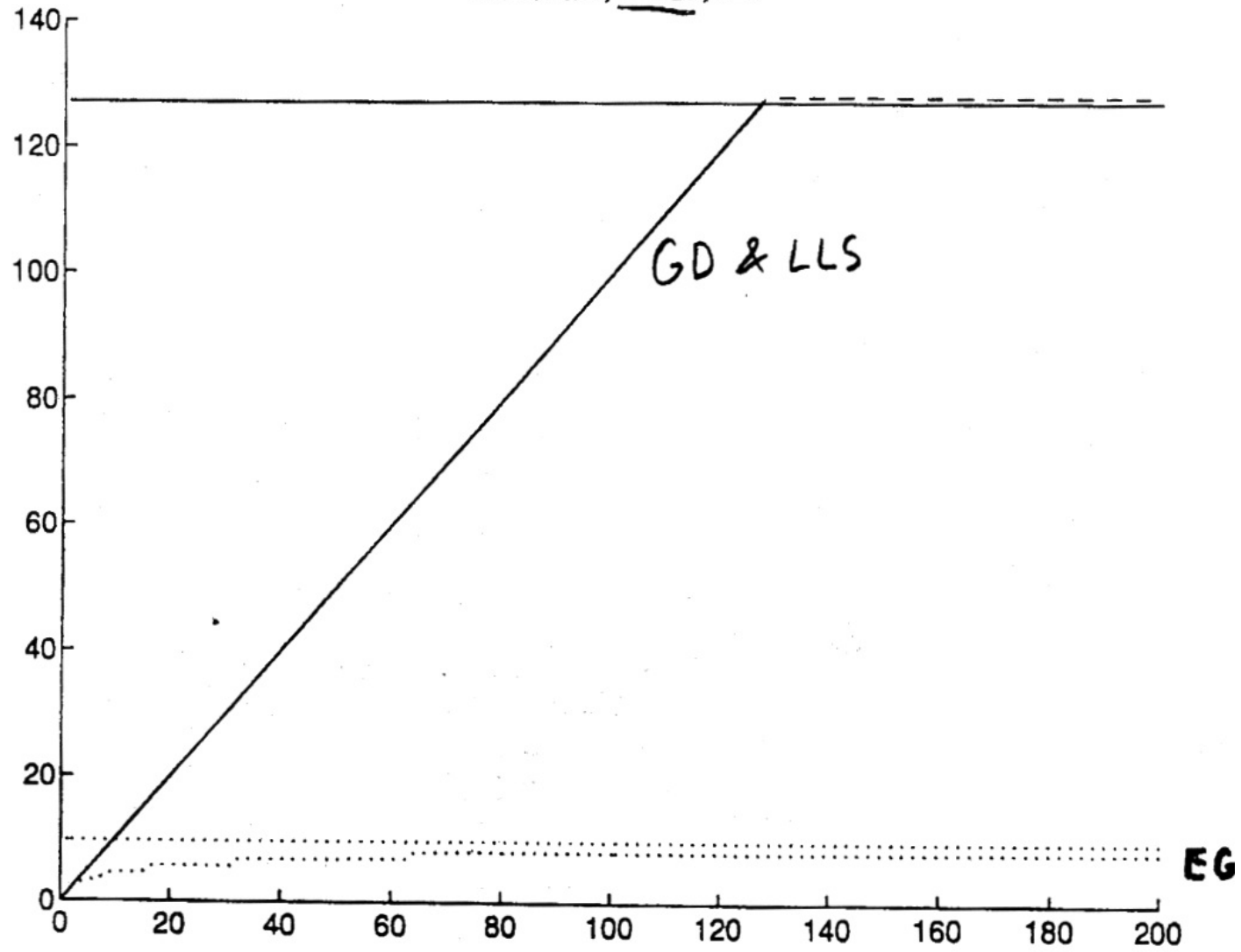
$$H_1 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \quad H_2 = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}, \quad \dots$$

$$H_{t+1} = \begin{pmatrix} H_t & H_t \\ H_t & -H_t \end{pmatrix}$$

All Rows are Orthogonal!

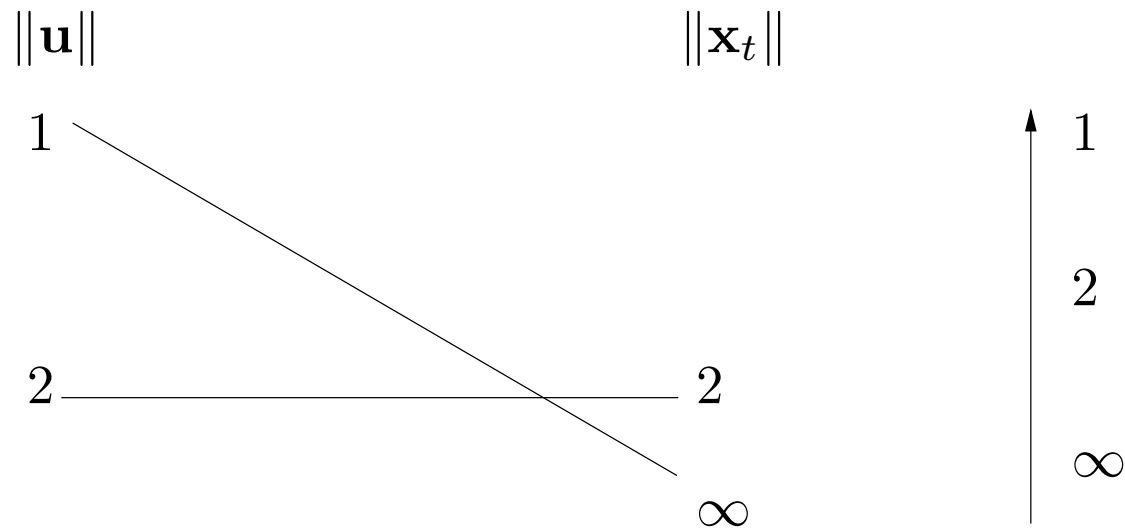
Use Rows as Instances!

Hadamard, $n=128$, $k=1$



Incomparable Loss bounds (cont)

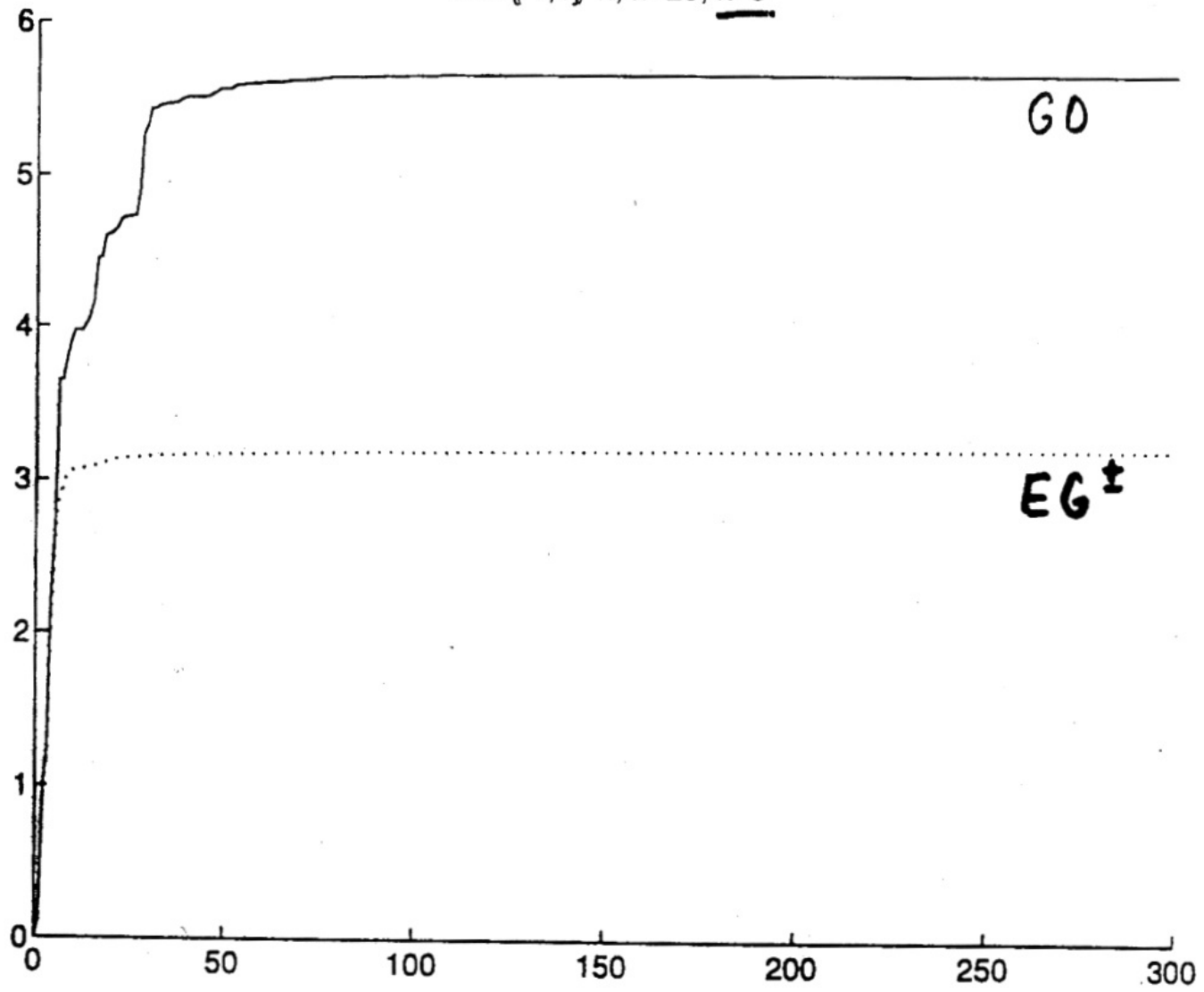
Gradient Descent has advantage:



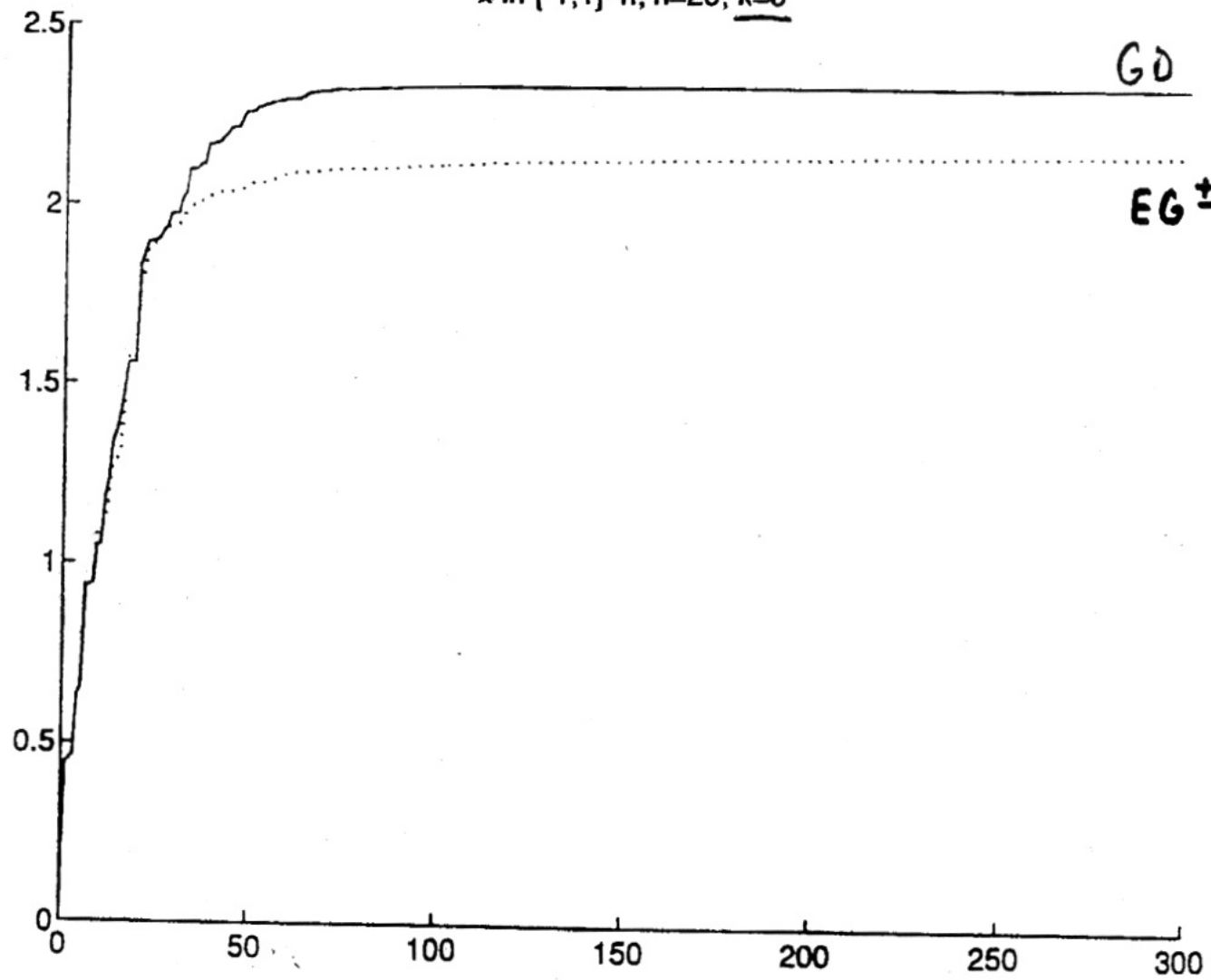
Widest gap between ℓ_1 and ℓ_2 norms if comparator \mathbf{u} is dense

...and there is the additional factor $\log N$ for EG

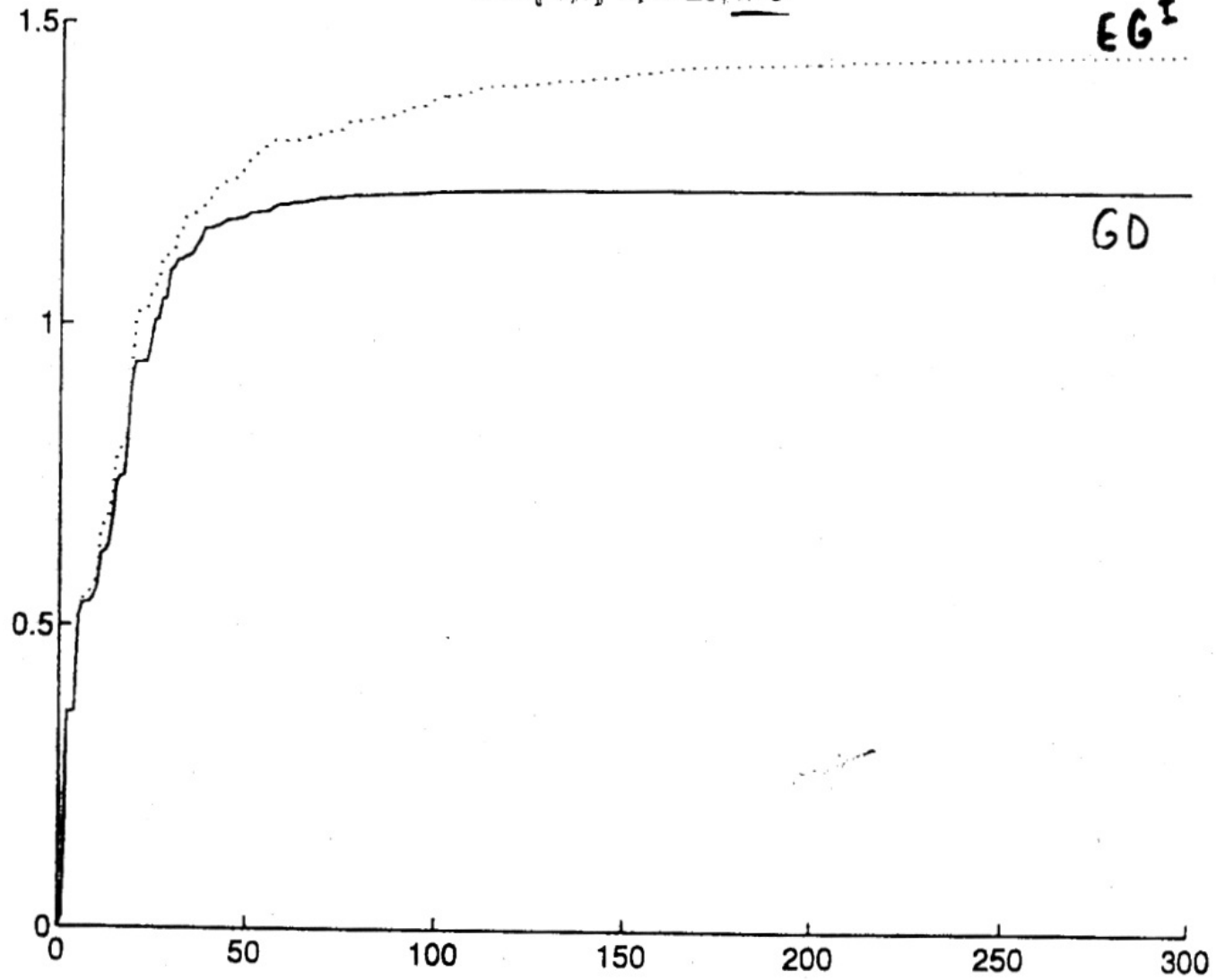
$x \in \{-1, 1\}^n; n=20; k=3$



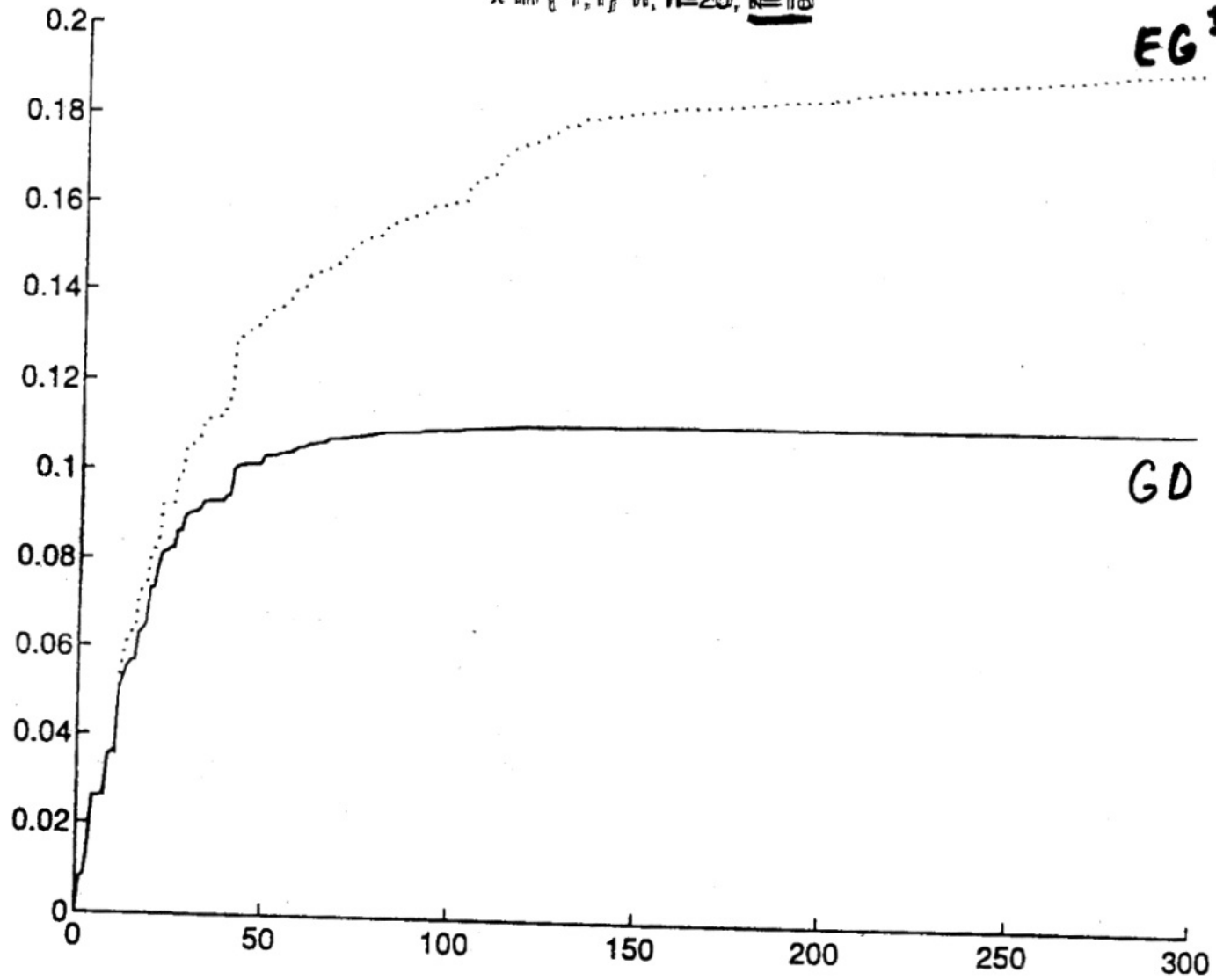
$x \text{ in } \{-1,1\}^n; n=20; \underline{k=6}$



$x \in \{-1, 1\}^n; n=20; k=9$



$x \text{ in } \{-1, 1\}^n; n=20; \underline{k=18}$



Summary of Comparison

- **EG** better when:
 - Instances \mathbf{x}_t are dense ($\|\mathbf{x}_t\|_\infty \ll \|\mathbf{x}_t\|_2$)
 - best weight vector is sparse ($\|\mathbf{u}_t\|_1 \approx \|\mathbf{u}_t\|_2$)
- **GD** better when:
 - instances are sparse ($\|\mathbf{x}_t\|_\infty \approx \|\mathbf{x}_t\|_2$)
 - best weight vector is dense ($\|\mathbf{u}_t\|_2 \ll \|\mathbf{u}_t\|_2$)

GD can be exponentially worse than EG