

A case where a spindly two-layer linear network  
~~whips~~ any neural network  
decisively outperforms  
with a fully connected input layer

Manfred K. Warmuth & Wojciech Kotłowski & Ehsan Amid  
Google and Poznan

ALT 2021  
March 16, 2021

## Paradigmatic sparse linear problem

$$\begin{pmatrix} -1 & -1 & 1 & -1 & -1 \\ 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & 1 & -1 & -1 \\ 1 & 1 & -1 & 1 & 1 \\ 1 & -1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$$

$\pm$  matrix random or Hadamard

After receiving example  $(\mathbf{x}_t, y_t)$   
and incurring loss  $1/2 (\mathbf{x}_t^\top \mathbf{w}_t - y_t)^2$ , update:

additive, GD:  $w_{t+1,i} = w_{t,i} - \underbrace{\eta \mathbf{x}_{t,i} (\mathbf{x}_t^\top \mathbf{w}_t - y_t)}_{\text{gradient}}$

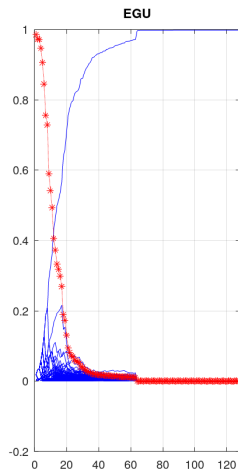
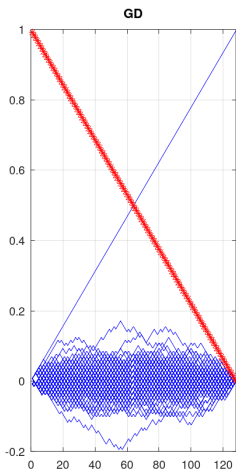
multiplicative, EGU:  $w_{t+1,i} = w_{t,i} \exp(-\eta \mathbf{x}_{t,i} (\mathbf{x}_t^\top \mathbf{w}_t - y_t))$

# Linear regression

Major differences in following paradigmatic setup:

**128x128 Hadamard matrix**

**Permuted** rows are instances, labels are any fixed column



x-axis:  $k = 1..128$

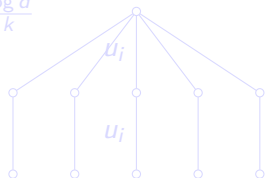
# Recent surprise and goal of this talk

Hadamard Problem: Instances are rows of the Hadamard matrix, labels are one of the features (e.g.  $\mathbf{w} = \mathbf{e}_1$ )

$$\begin{pmatrix} + & + & + & + \\ + & - & + & - \\ + & + & - & - \\ + & - & - & + \end{pmatrix}$$

**Conjecture** [DW14]: For the Hadamard problem all neural nets trained w. GD incur loss at least  $1 - k/d$  after seeing  $k$  examples

**Surprise:** Spindly GD trained linear net cracks shifted Hadamard with loss  $\leq \frac{\log d}{k}$



Simulates EGU [AW20]

Square reparameterization trick reintroduced in [GWBNS17]

**Here:** Lower bound of  $\geq 1 - k/d$  for any GD trained net with a fully connected input layer

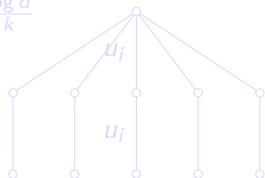
# Recent surprise and goal of this talk

Hadamard Problem: Instances are rows of the Hadamard matrix, labels are one of the features (e.g.  $\mathbf{w} = \mathbf{e}_1$ )

$$\begin{pmatrix} + & + & + & + \\ + & - & + & - \\ + & + & - & - \\ + & - & - & + \end{pmatrix}$$

**Conjecture [DW14]:** For the Hadamard problem all neural nets trained w. GD incur loss at least  $1 - k/d$  after seeing  $k$  examples

**Surprise:** Spindly GD trained linear net cracks shifted Hadamard with loss  $\leq \frac{\log d}{k}$



Simulates EGU [AW20]

Square reparameterization trick reintroduced in [GWBNS17]

**Here:** Lower bound of  $\geq 1 - k/d$  for any GD trained net with a fully connected input layer

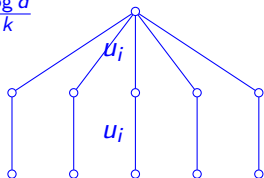
# Recent surprise and goal of this talk

Hadamard Problem: Instances are rows of the Hadamard matrix, labels are one of the features (e.g.  $\mathbf{w} = \mathbf{e}_1$ )

$$\begin{pmatrix} + & + & + & + \\ + & - & + & - \\ + & + & - & - \\ + & - & - & + \end{pmatrix}$$

**Conjecture** [DW14]: For the Hadamard problem all neural nets trained w. GD incur loss at least  $1 - k/d$  after seeing  $k$  examples

**Surprise:** Spindly GD trained linear net cracks shifted Hadamard with loss  $\leq \frac{\log d}{k}$



Simulates EGU [AW20]

Square reparameterization trick reintroduced in [GWBNS17]

**Here:** Lower bound of  $\geq 1 - k/d$  for any GD trained net with a fully connected input layer

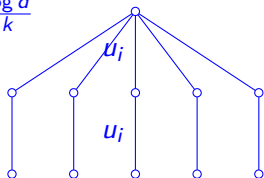
# Recent surprise and goal of this talk

Hadamard Problem: Instances are rows of the Hadamard matrix, labels are one of the features (e.g.  $\mathbf{w} = \mathbf{e}_1$ )

$$\begin{pmatrix} + & + & + & + \\ + & - & + & - \\ + & + & - & - \\ + & - & - & + \end{pmatrix}$$

**Conjecture [DW14]:** For the Hadamard problem all neural nets trained w. GD incur loss at least  $1 - k/d$  after seeing  $k$  examples

**Surprise:** Spindly GD trained linear net cracks shifted Hadamard with loss  $\leq \frac{\log d}{k}$



Simulates EGU [AW20]

Square reparameterization trick reintroduced in [GWBNS17]

**Here:** Lower bound of  $\geq 1 - k/d$  for any GD trained net with a **fully connected input layer**

# Old versus new lower bounds

Previous SVD based lower bound technique of [WV05]

- + Holds for any expansion
- Only for single linear neuron
- Can't go beyond 2 layer linear
- Restricted to square loss

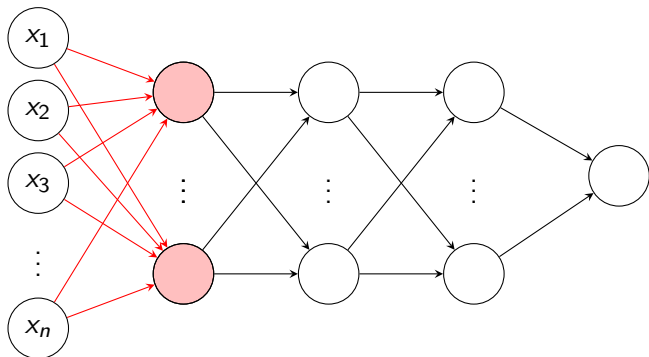
NEW: Hold for any GD trained any neural net **of any depth** with a fully connected input layer

- Does not hold for any expansion (essentially Hadamard or Gaussian)
- + Holds for single target such as constant feature  $e_1$
- + Very general losses

**Crux: GD trained w. fully connected input layer implies rotational invariance of the predictions**



input layer 0    hidden layer 1    hidden layer 2    hidden layer 3    output layer 4



- fully connected input layer with rotation invariant initialization at 1st hidden layer
- weights of first hidden layer trained with GD, any learning rates
- one output node, any initialization
- otherwise any architecture and initialization
- any differentiable transfer functions at the internal nodes
- very general loss function

# Minimal loss function assumptions

$L(y, \hat{y})$  differentiable in  $\hat{y}$  and  
 $\min_{\hat{y}} \frac{L(-1, \hat{y}) + L(+1, \hat{y})}{2}$  is some positive constant  $c$

Simplifying assumption for the talk: loss is **convex**

Then for proving lower bounds, any randomized algorithm  $\hat{y}$  can be turned into a **deterministic algorithm**  $\hat{y}_{\text{det}}$ :

$$\hat{y}_{\text{det}}(\mathbf{x} | (\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}})) = \mathbb{E} [\hat{y}(\mathbf{x} | (\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}))]$$

which by Jensen's inequality has loss no greater than the expected loss of  $\hat{y}$  on any instance

For simplicity we use **square loss**, i.e.  $c = 1$

# It's all about rotation invariance of the predictions

A prediction algorithm  $\hat{y}(\mathbf{x} | (\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}))$  is called **rotation invariant** if

$$\hat{y}(\mathbf{U}\mathbf{x} | (\mathbf{X}_{\text{tr}}\mathbf{U}^{\text{T}}, \mathbf{y}_{\text{tr}})) = \hat{y}(\mathbf{x} | (\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}})), \text{ for any orthogonal matrix } \mathbf{U}$$

GD trained neural nets with fully connected input layers are rotation invariant because gradients/weights at input nodes are linear combinations of the training instances

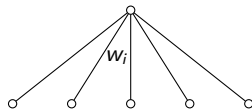
$$\frac{\partial f(\mathbf{x} \cdot \mathbf{w})}{\partial \mathbf{w}} = f'(a)|_{a=\mathbf{x} \cdot \mathbf{w}} \mathbf{x}$$

**Crux:** Any input node that receives an input feature, must receive all input features

Newton, versions of AdaGrad & Adam, ..., also rotation invariant

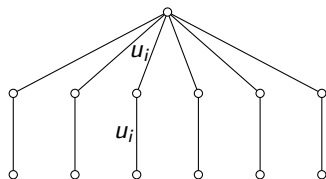
# Interlude

EGU on single neuron



$$w_{t+1,i} = w_{t,i} \exp\left(-\eta \mathbf{x}_{t,i}(\mathbf{x}_t^\top \mathbf{w}_t - y_t)\right)$$

and GD on spindly network



$$u_{t+1,i}^2 = u_{t,i}^2 \left(1 - \eta/2 \mathbf{x}_{t,i}(\mathbf{x}_t^\top \mathbf{w}_t - y_t)\right)^2$$

both decidedly not rotation invariant:

- ▶ access individual features
- ▶ rank grows exponentially fast

# Hardness for any rotation invariant algorithm

Initial problem  $(\mathbf{H}, \mathbf{1})$ .

Too easy

For each sign pattern  $\mathbf{s} \in \{+1, -1\}^d$  consider problem  $(\text{diag}(\mathbf{s})\mathbf{H}, \mathbf{s})$

Define a rotation matrix  $\mathbf{U}_s$  as  $1/\sqrt{d} \text{diag}(\mathbf{s})\mathbf{H}$

The predictions of any rotation invariant algorithm on

$(\text{diag}(\mathbf{s})\mathbf{H}, \mathbf{s})$  and  $(\text{diag}(\mathbf{s})\mathbf{H}\mathbf{U}_s^T, \mathbf{s}) = (\sqrt{d}\mathbf{I}, \mathbf{s})$  are the same

Now fix  $\mathbf{s}_{1:k}$ . The algorithm receives the same first  $k$  training examples for each problem  $(\sqrt{d}\mathbf{I}, \mathbf{s})$ . Also since  $\mathbf{s}_{(k+1):d}$  is chosen uniformly, each of the  $d - k$  unseen examples is labeled  $\pm 1$  with equal probability. So the best prediction on these  $d - k$  examples is 0, incurring square loss at least 1 for each unseen example

Conclusion: Expected average loss on all  $d$  examples is at least  $(d - k)/d = 1 - k/d$

# Hardness for any rotation invariant algorithm

Initial problem  $(\mathbf{H}, \mathbf{1})$ .  
 $\substack{d,d \\ d,1}$

Too easy

For each sign pattern  $\mathbf{s} \in \{+1, -1\}^d$  consider problem  $(\text{diag}(\mathbf{s})\mathbf{H}, \mathbf{s})$

Define a rotation matrix  $\mathbf{U}_s$  as  $1/\sqrt{d} \text{diag}(\mathbf{s})\mathbf{H}$

The predictions of any rotation invariant algorithm on

$(\text{diag}(\mathbf{s})\mathbf{H}, \mathbf{s})$  and  $(\text{diag}(\mathbf{s})\mathbf{H}\mathbf{U}_s^\top, \mathbf{s}) = (\sqrt{d}\mathbf{I}, \mathbf{s})$  are the same

Now fix  $\mathbf{s}_{1:k}$ . The algorithm receives the same first  $k$  training examples for each problem  $(\sqrt{d}\mathbf{I}, \mathbf{s})$ . Also since  $\mathbf{s}_{(k+1):d}$  is chosen uniformly, each of the  $d - k$  unseen examples is labeled  $\pm 1$  with equal probability. So the best prediction on these  $d - k$  examples is 0, incurring square loss at least 1 for each unseen example

Conclusion: Expected average loss on all  $d$  examples is at least  $(d - k)/d = 1 - k/d$

- ▶ Lower bound holds even for rotation invariant initialization of hidden nodes of first layer  
(which must all be connected to all inputs)
- ▶ There are versions of the problem where the features are 0/1 instead of  $\pm 1$
- ▶ Only minimal requirements on the loss needed:  
 $L(y, \hat{y})$  differentiable in  $\hat{y}$  and  
 $\min_{\hat{y}} \frac{L(-1, \hat{y}) + L(+1, \hat{y})}{2}$  is some positive constant  $c$

## Extensions cont.

- ▶ Only slightly weaker linear lower bound for the case when entries of  $\mathbf{X}$  are i.i.d. Gaussian and the target is first column  $\mathbf{X}\mathbf{e}_1$ :

After seeing  $k$  examples, the expected average square loss on all  $d$  examples is at least  $(1 - k/d)^2 = 1 - 2k/d + (k/d)^2$   
**(so far only for square loss)**

- ▶ Experimentally the same lower bound for square loss holds for random  $\pm 1$  matrices and single feature targets, **but no proof yet**



# Is there a method here?

Structure of neural net / update

→ invariance

→ lower bound

	additive	multiplicative
rot. invariance	$\hat{y}(\mathbf{U}\mathbf{x}   (\mathbf{X}_{\text{tr}} \mathbf{U}^{\top}, \mathbf{y}_{\text{tr}}))$ $= \hat{y}(\mathbf{x}   (\mathbf{X}_{\text{tr}}, \mathbf{y}_{\text{tr}}))$	$\hat{y}(\mathbf{U}\mathbf{x}\mathbf{z}^{\top} \mathbf{V}^{\top}   (\mathbf{U}\mathbf{X}_{\text{tr}} \mathbf{Z}_{\text{tr}}^{\top} \mathbf{V}^{\top}, \mathbf{y}_{\text{tr}}))$ $= \hat{y}(\mathbf{x}\mathbf{z}^{\top}   (\mathbf{X}_{\text{tr}} \mathbf{Z}_{\text{tr}}^{\top}, \mathbf{y}_{\text{tr}}))$
linear comb.	$\mathbf{w} = \mathbf{X}_{\text{tr}}^{\top} \mathbf{a}$	$\mathbf{W} = \mathbf{X}_{\text{tr}} \mathbf{C} \mathbf{Z}_{\text{tr}}^{\top}$
hard problem	Hadamard	???
lower bound	$1 - k/d$	???

[WKZ14]

# The future is open

For Hadamard instances:

- ▶ We showed above that GD trained neural nets with a complete input layer cannot learn a **single target** sample efficiently.
- ▶ However **any single target feature  $\mathbf{y}$  can** be learned when inputs are transformed by the map  $\phi(\mathbf{H}) = \mathbf{y}$
- ▶ **Conjecture: The  $d$  target features cannot be learned with any  $\phi(\cdot)$  map when we average over targets**

## Rotation invariance of on-line Newton

**Inductive hypothesis:** Under transformation  $\mathbf{x}_t \mapsto \mathbf{U}\mathbf{x}_t$  for all  $t$  we have  $\mathbf{w}_t \mapsto \mathbf{U}\mathbf{w}_t$ .

Suppose it holds for  $\mathbf{w}_j, j \leq t$  and we show it holds for  $\mathbf{w}_{t+1}$ .

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{A}_t^{-1} \nabla_t, \quad \mathbf{A}_t = \epsilon \mathbf{I} + \sum_{j \leq t} \nabla_j \nabla_j^\top, \quad \nabla_t = (\hat{y}_t - y_t) \mathbf{x}_t$$

Since  $\hat{y}_j = \mathbf{w}_j^\top \mathbf{x}_j \mapsto \mathbf{w}_j^\top \mathbf{U}^\top \mathbf{U} \mathbf{x}_j = \hat{y}_j$ , we have  $\nabla_j \mapsto \mathbf{U} \nabla_j$  for  $j \leq t$ , so  $\mathbf{A}_t \mapsto \mathbf{U} \mathbf{A}_t \mathbf{U}^\top$ . Therefore:

$$\mathbf{w}_{t+1} \mapsto \mathbf{U} \mathbf{w}_t - \eta \mathbf{U} \mathbf{A}_t^{-1} \mathbf{U}^\top \mathbf{U} \nabla_t = \mathbf{U} \mathbf{w}_{t+1}$$

- ▶ Diagonal version:

$$\mathbf{w}_{t+1,i} = \mathbf{w}_{t,i} - \frac{\eta}{\sqrt{\epsilon + \sum_{j \leq t} \nabla_{j,i}^2}} \nabla_{t,i}$$

Clearly not rotation-invariant, but when  $\mathbf{x}_{t,i} = \pm 1$  for all  $t, i$ , then  $\nabla_{j,i}^2 = (y_t - \hat{y}_t)^2$ , i.e. the effective learning rate is shared among all coordinates and rotation invariance holds.

- ▶ Full version:

$$\mathbf{w}_t = \mathbf{w}_t - \eta \sqrt{\mathbf{A}_t} \nabla_t, \quad \mathbf{A}_t = \epsilon \mathbf{I} + \sum_{j \leq t} \nabla_j \nabla_j^\top$$

Rotation invariant, proof analogous to On-line Newton

- ▶ Adagrad with momentum and exponentially decaying memory on the past
- ▶ Essentially the same arguments as for diagonal Adagrad: generally not rotation invariant, but for  $\pm 1$  valued data it is rotation invariant.

$$\mathbf{w}_t = \mathbf{A}_t^{-1} \sum_{j \leq t-1} y_j \mathbf{x}_j, \quad \mathbf{A}_t = \epsilon \mathbf{I} + \sum_{j \leq t} \mathbf{x}_j \mathbf{x}_j^\top.$$

VAW is also rotation invariant: under transformation  $\mathbf{x}_t \mapsto \mathbf{U} \mathbf{x}_t$  for all  $t$ , we clearly have  $\mathbf{A}_t \mapsto \mathbf{U} \mathbf{A}_t \mathbf{U}^\top$ , so  $\mathbf{A}_t^{-1} \mapsto \mathbf{U} \mathbf{A}_t^{-1} \mathbf{U}^\top$ , and thus

$$\mathbf{w}_t \mapsto \mathbf{U} \mathbf{A}_t^{-1} \mathbf{U}^\top \sum_{j \leq t-1} y_j \mathbf{U} \mathbf{x}_j = \mathbf{U} \mathbf{A}_t^{-1} \sum_{j \leq t-1} y_j \mathbf{x}_j = \mathbf{U} \mathbf{w}_t$$