

Open Problem:
Learning sparse linear concepts
by priming the features

Manfred K. Warmuth and Ehsan Amid
Google Inc.

July 15, 2023
COLT 2023 in Bangalore

Multiplicative & additive updates

Single example (\mathbf{x}, y) for linear regression

Square loss $(\mathbf{x} \cdot \mathbf{w} - y)^2/2$

$$w_i = w_i - \eta(\mathbf{w} \cdot \mathbf{x} - y)x_i \quad (\text{GD})$$

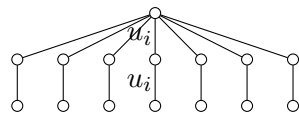
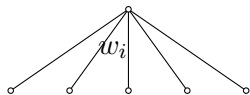
$$w_i = w_i \exp^{-\eta(\mathbf{w} \cdot \mathbf{x} - y)x_i} \quad (\text{EGUnnormalized})$$

$$w_i = \frac{w_i \exp^{-\eta(\mathbf{w} \cdot \mathbf{x} - y)x_i}}{\sum_j w_j \exp^{-\eta(\mathbf{w} \cdot \mathbf{x} - y)x_j}} \quad (\text{EG})$$

Key: Multiplicative update much much better at learning when label y is close to sparse linear

$\log d$ versus d dependence in the regret

EGU approximated by GD on spindly network



$$w_i = w_i \exp^{-\eta(\mathbf{w} \cdot \mathbf{x} - y)x_i}$$

\approx

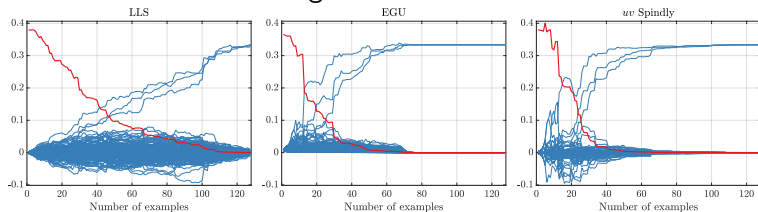
$$u_i = u_i - \eta(\mathbf{u} \circ \mathbf{u} \cdot \mathbf{x} - y)x_i u_i$$

EGU

$\mathbf{u} \circ \mathbf{u}$ spindly

$$\dot{\log}(\mathbf{w}) = -\eta(\mathbf{w} \cdot \mathbf{x} - y)\mathbf{x} \text{ equivalent } \dot{\mathbf{u}} = -(\mathbf{u} \circ \mathbf{u} \cdot \mathbf{x} - y)\mathbf{u} \circ \mathbf{x}$$

In many cases same regret bounds

Average of $3 \pm$ features

We draw a random $n \times n \pm 1$ matrix \mathbf{X} , with $n = 128$. The rows are the n instances and the sparse label vector \mathbf{y} is the average of the first three column features. After seeing $t = 1 : 128$ examples, we plot weights (in blue) after applying LLS or the priming methods to all examples. For the online update, we do multiple passes over the past examples until the weight vector is consistent (EGU \pm is a 2-sided version of EGU). In red, we always plot the average square loss on all 128 examples. LLS picks up the target weights too slowly.

It is all about rotation invariance

Any neural net with a complete input layer and rotation invariant initialization can't learn sparse linear efficiently
when trained with GD

So on sparse targets, beaten by GD on spindly

One big advantage of the GD family

Closed form in batch case:

(\mathbf{X}, \mathbf{y}) , instances row-wise

$$\mathbf{w}_{\text{LLS}} = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{X}\mathbf{w} - \mathbf{y})^2 = \mathbf{X}^\dagger \mathbf{y} \quad (\text{LLS})$$

Is there a closed form method that approximates multiplicative updates, i.e. is good at sparse targets

i th feature multiplied by the prime factor p_i & apply LLS:

$$\mathbf{w}_p = \text{diag}(\mathbf{p}) (\mathbf{X} \text{diag}(\mathbf{p}))^\dagger \mathbf{y}$$

Choices

1. $p_i = \text{argmin}_{p_i} (\mathbf{X}(:, i) p_i - \mathbf{y})^2 = X(:, i)^\dagger \mathbf{y} = \frac{X(:, i)^\top}{\|X(:, i)\|^2} \mathbf{y}$
(1-dim LLS per i th feature)
2. $p_i = \frac{(X(:, i) - \overline{X(:, i)}) (y - \bar{y})}{\sqrt{(X(:, i) - \overline{X(:, i)})^2} \sqrt{y - \bar{y}}}$ (Pearson correlation coefficients)
3. $\mathbf{p} = \mathbf{w}_{\text{LLS}} = \mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y}$ (LLS for priming as well)

i th feature multiplied by the prime factor p_i & apply LLS:

$$\mathbf{w}_p = \text{diag}(\mathbf{p}) (\mathbf{X} \text{diag}(\mathbf{p}))^\dagger \mathbf{y}$$

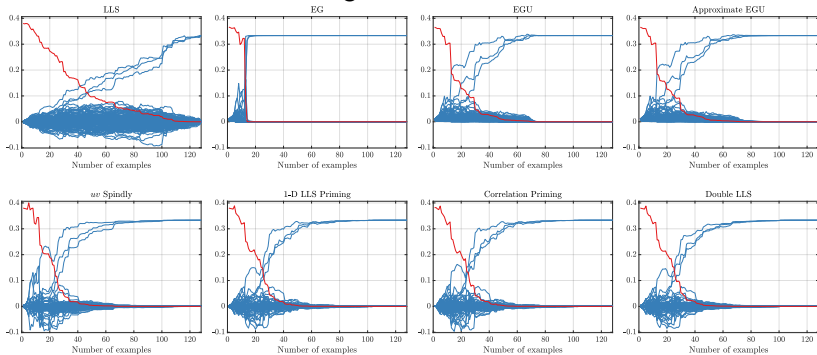
Choices

1. $p_i = \text{argmin}_{p_i} (\mathbf{X}(:, i) p_i - \mathbf{y})^2 = X(:, i)^\dagger \mathbf{y} = \frac{X(:, i)^\top \mathbf{y}}{\|X(:, i)\|^2}$
(1-dim LLS per i th feature)
2. $p_i = \frac{(X(:, i) - \overline{X(:, i)}) (y - \overline{y})}{\sqrt{(X(:, i) - \overline{X(:, i)})^2} \sqrt{y - \overline{y}}}$ (Pearson correlation coefficients)
3. $\mathbf{p} = \mathbf{w}_{\text{LLS}} = \mathbf{X}^\dagger \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^\dagger \mathbf{X}^\top \mathbf{y}$ (LLS for priming as well)

Open

- Are there competitive regret bound for any of the priming methods?
- What is the optimal priming function for sparse linear problems?
- Are there priming methods for learning sparse disjunctions?

Average of 3 \pm features



LLS picks up the target weights too slowly. All three priming methods perform similarly to the multiplicative updates

It does not help for learning sparse concepts

Open

Can kernels be primed efficiently?

Dream

If priming has regret bounds that grow with $\log d$
then learn sparse concepts in feature space

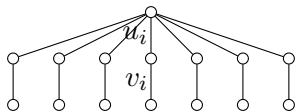
The Hoff $\mathbf{u} \circ \mathbf{v}$ Lemma

$$\begin{aligned} & \underset{\mathbf{w}}{\operatorname{argmin}} h(\mathbf{w}) + 1/\eta \|\mathbf{w}\|_1 \\ & = \underset{\mathbf{u}, \mathbf{v}}{\operatorname{argmin}} h(\mathbf{u} \circ \mathbf{v}) + 1/\eta (1/2\mathbf{u}^2 + 1/2\mathbf{v}^2) \end{aligned}$$

- A one iteration batch update
- Assumes initializations are $\mathbf{w}, \mathbf{u}, \mathbf{v} = \mathbf{0}$

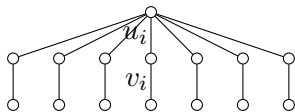
Big picture

- Multiplicative = spindly in continuous case
- Spindly updates seem to have an L_1 motivation as well
- They are regret bounds for some spindlified updates
- Priming can be viewed as layerwise spindly



Big picture

- Multiplicative = spindly in continuous case
- Spindly updates seem to have an L_1 motivation as well
- They are regret bounds for some spindlified updates
- Priming can be viewed as layerwise spindly



**How many hacks can we get away with
and still learn sparse concepts?**