# Follow the leader with Dropout perturbations - Additive versus multiplicative noise

Manfred K. Warmuth

UCSC

Simon's Institute
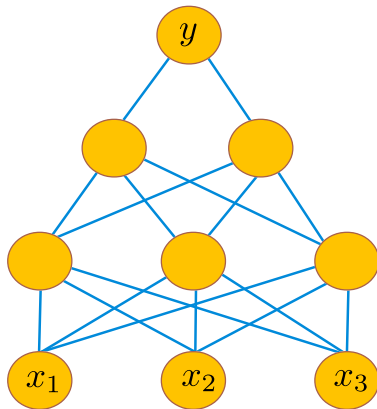Nov. 17, 2016

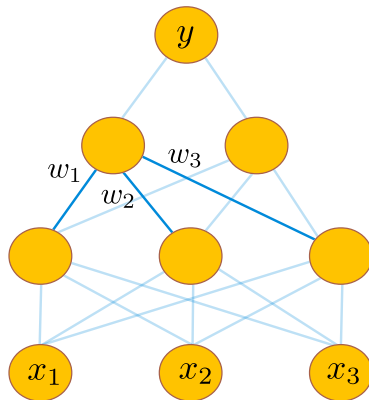Joint work with    Tim Van Erven    and    Wojciech Kotłowski

Universiteit Leiden

Major insights from [Devroye, Lugosi, Neu 2013]
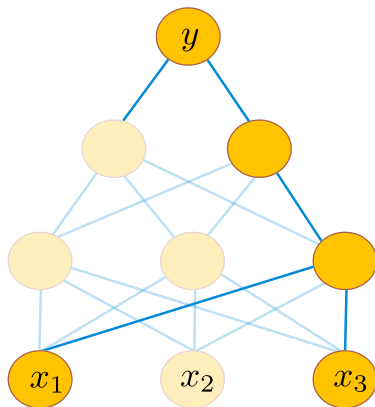
# Dropout training



- Stochastic gradient descent
- Randomly remove every hidden/input node with prob. $\frac{1}{2}$ before each gradient descent update

[Hinton et al. 2012]

# Dropout training



- Very successful in image recognition & speech recognition
- Why does it work?

  [Wagner, Wang, Liang 2013]
  [Helmbold, Long 2014]

Prove bounds for dropout
- single neuron
- linear loss

# Online learning with expert

|        | $E_1$ | $E_2$ | $E_3$ | $\ldots$ | $E_n$ | $prediction$ | label | loss |
|--------|-------|-------|-------|----------|-------|--------------|-------|------|
| day 1  | 0     | 1     | 0     | $\ldots$ | 0     | 0            | 1     | 1    |

## Online learning with expert

|       | $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ | $prediction$ | label | loss |
|-------|-------|-------|-------|-----|-------|--------------|-------|------|
| day 1 | 0     | 1     | 0     | ... | 0     | 0            | 1     | 1    |
| day 2 | 1     | 1     | 0     | ... | 0     | 1            | 1     | 0    |

# Online learning with expert

| | $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ | $prediction$ | label | loss |
|---|---|---|---|---|---|---|---|---|
| day 1 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 1 |
| day 2 | 1 | 1 | 0 | ... | 0 | 1 | 1 | 0 |
| | | | | | | | | |
| notation | $x_1$ | $x_1$ | $x_2$ | ... | $x_n$ | $\widehat{y}$ | $y$ | $|\widehat{y} - y|$ |

# Online learning with expert

| | $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ | $prediction$ | label | loss |
|---|---|---|---|---|---|---|---|---|
| day 1 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 1 |
| day 2 | 1 | 1 | 0 | ... | 0 | 1 | 1 | 0 |
| | | | | | | | | |
| notation | $x_1$ | $x_1$ | $x_2$ | ... | $x_n$ | $\widehat{y}$ | $y$ | $|\widehat{y} - y|$ |
| scope | $\in [0,1]$ | | | ... | | $\in [0,1]$ | $\in \{0,1\}$ | $\in [0,1]$ |

# Online learning with expert

| | $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ | *prediction* | label | loss |
|---|---|---|---|---|---|---|---|---|
| day 1 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 1 |
| day 2 | 1 | 1 | 0 | ... | 0 | 1 | 1 | 0 |
| | | | | | | | | |
| notation | $x_1$ | $x_1$ | $x_2$ | ... | $x_n$ | $\widehat{y}$ | $y$ | $|\widehat{y} - y|$ |
| scope | $\in [0,1]$ | | | ... | | $\in [0,1]$ | $\in \{0,1\}$ | $\in [0,1]$ |

- Algorithm maintains probability vector $\boldsymbol{w}$:
  - prediction $\widehat{y} = \boldsymbol{w} \cdot \boldsymbol{x}$

# Online learning with expert

| | $E_1$ | $E_2$ | $E_3$ | ... | $E_n$ | $prediction$ | label | loss |
|---|---|---|---|---|---|---|---|---|
| day 1 | 0 | 1 | 0 | ... | 0 | 0 | 1 | 1 |
| day 2 | 1 | 1 | 0 | ... | 0 | 1 | 1 | 0 |
| | | | | | | | | |
| notation | $x_1$ | $x_1$ | $x_2$ | ... | $x_n$ | $\widehat{y}$ | $y$ | $\|\widehat{y} - y\|$ |
| scope | $\in [0,1]$ | | | ... | | $\in [0,1]$ | $\in \{0,1\}$ | $\in [0,1]$ |

- Algorithm maintains probability vector $\boldsymbol{w}$:
  - prediction $\widehat{y} = \boldsymbol{w} \cdot \boldsymbol{x}$

- Loss linear because label $y \in \{0,1\}$

- $\underbrace{|\overbrace{\boldsymbol{w} \cdot \boldsymbol{x}}^{\widehat{y}} - y|}_{\text{loss of alg.}} = \sum_i w_i \underbrace{|x_i - y|}_{\text{loss of expert } i}$
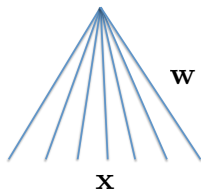
**Predicting with expert advice**

$$\hat{y} = \mathbf{w} \cdot \mathbf{x} \qquad \text{loss } |\hat{y} - y|$$

**Predicting with expert advice**

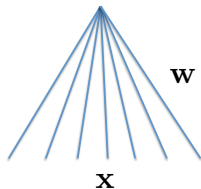$$\hat{y} = \mathbf{w} \cdot \mathbf{x} \qquad \text{loss } |\hat{y} - y|$$



trial $t$
- get advice vector $\boldsymbol{x}_t$
- predict $\widehat{y}_t = \boldsymbol{w}_t \cdot \boldsymbol{x}_t$
- get label $y_t$
- exp. losses $|x_{t,i} - y_t|$
- alg. loss $|\widehat{y}_t - y_t|$
- update $\boldsymbol{w}_t \to \boldsymbol{w}_{t+1}$
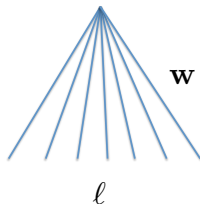
# On-line learning

**Predicting with expert advice**

$\hat{y} = \mathbf{w} \cdot \mathbf{x}$       $\text{loss } |\hat{y} - y|$



**Hedge setting**

$\text{loss } \mathbf{w} \cdot \ell$



trial $t$
- get advice vector $\boldsymbol{x}_t$
- predict $\widehat{y}_t = \boldsymbol{w}_t \cdot \boldsymbol{x}_t$
- get label $y_t$
- exp. losses $|x_{t,i} - y_t|$
- alg. loss $|\widehat{y}_t - y_t|$
- update $\boldsymbol{w}_t \rightarrow \boldsymbol{w}_{t+1}$

# On-line learning

**Predicting with expert advice**

$$\hat{y} = \mathbf{w} \cdot \mathbf{x} \qquad \text{loss } |\hat{y} - y|$$



$\mathbf{w}$

$\mathbf{x}$

trial $t$
- get advice vector $\boldsymbol{x}_t$
- predict $\widehat{y}_t = \boldsymbol{w}_t \cdot \boldsymbol{x}_t$
- get label $y_t$
- exp. losses $|x_{t,i} - y_t|$
- alg. loss $|\widehat{y}_t - y_t|$
- update $\boldsymbol{w}_t \to \boldsymbol{w}_{t+1}$

**Hedge setting**

$$\text{loss } \mathbf{w} \cdot \boldsymbol{\ell}$$



$\mathbf{w}$

$\boldsymbol{\ell}$

trial $t$
-
- predict $\boldsymbol{w}_t$
- get loss vector $\boldsymbol{\ell}_t$
- exp. losses $\ell_{t,i}$
- alg. loss $\boldsymbol{w}_t \cdot \boldsymbol{\ell}_t$
- update $\boldsymbol{w}_t \to \boldsymbol{w}_{t+1}$

trial $t$

- predict $\boldsymbol{w}_t$          or predict with random expert $\widehat{i}_t$

trial $t$
- predict $\boldsymbol{w}_t$        or predict with random expert $\widehat{i}_t$
- get loss vector $\boldsymbol{\ell}_t$
- alg. loss $\boldsymbol{w}_t \cdot \boldsymbol{\ell}_t$      or alg. expected loss $\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t} \cdot \boldsymbol{\ell}_t\right] = \underbrace{\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t}\right]}_{\boldsymbol{w}_t} \cdot \boldsymbol{\ell}_t$

trial $t$

- predict $\boldsymbol{w}_t$            or predict with random expert $\widehat{i}_t$
- get loss vector $\boldsymbol{\ell}_t$

- alg. loss $\boldsymbol{w}_t \cdot \boldsymbol{\ell}_t$      or alg. expected loss $\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t} \cdot \boldsymbol{\ell}_t\right] = \underbrace{\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t}\right]}_{\boldsymbol{w}_t} \cdot \boldsymbol{\ell}_t$

- update $\boldsymbol{w}_t \rightarrow \boldsymbol{w}_{t+1}$

trial $t$

- predict $\boldsymbol{w}_t$          or predict with random expert $\widehat{i}_t$
- get loss vector $\boldsymbol{\ell}_t$
- alg. loss $\boldsymbol{w}_t \cdot \boldsymbol{\ell}_t$     or alg. expected loss $\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t} \cdot \boldsymbol{\ell}_t\right] = \underbrace{\mathbb{E}\left[\boldsymbol{e}_{\widehat{i}_t}\right]}_{\boldsymbol{w}_t} \cdot \boldsymbol{\ell}_t$

- update $\boldsymbol{w}_t \rightarrow \boldsymbol{w}_{t+1}$

                weights are implicit

Only works for linear loss

Worst-case **regret**

$$\underbrace{\sum_{t=1}^{T} \boldsymbol{w}_t \cdot \boldsymbol{\ell}_t}_{\text{total expected loss of alg}} \quad - \quad \underbrace{\inf_i \ell_{\leq T, i}}_{\text{loss } \ell^* \text{ of best expert}}$$

Should be logarithmic in $\#$ of experts $n$

# Main algorithms

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | 1 | 0 | 0 | 1 |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | 1 | 1 | 1 |

$\ell_{\leq t-1,i}$

# Main algorithms

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | 1 | 0 | 0 | 1 |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | 1 | 1 | 1 |
| $\ell_{\leq t-1,i}$ | 1 | 2 | 1 | 2 | 3 |

# Main algorithms

|           | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|-----------|-------|-------|-------|-------|-------|
|           | 0     | 1     | 0     | 0     | 1     |
|           | 1     | 1     | 0     | 1     | 1     |
| day $t-1$ | 0     | 0     | 1     | 1     | 1     |
| $\ell_{\leq t-1,i}$ | 1 | 2 | 1 | 2 | 3 |

FL $\qquad \widehat{i}_t = \operatorname{argmin}_i \ \ell_{\leq t-1,i}$ $\qquad$ ties broken uniformly

# Main algorithms

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | 1 | 0 | 0 | 1 |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | 1 | 1 | 1 |
| $\ell_{\leq t-1,i}$ | 1 | 2 | 1 | 2 | 3 |

FL $\qquad \widehat{i}_t = \operatorname{argmin}_i \, \ell_{\leq t-1,i}$ $\qquad\qquad$ ties broken uniformly

FPL($\eta$) $\quad \widehat{i}_t = \operatorname{argmin}_i \, \ell_{\leq t-1,i} + \frac{1}{\eta}\xi_{t,i}$ $\quad$ indep. <u>additive</u> noise

# Main algorithms

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | 1 | 0 | 0 | 1 |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | 1 | 1 | 1 |
| $\ell_{\leq t-1,i}$ | 1 | 2 | 1 | 2 | 3 |

FL $\qquad \widehat{i}_t = \text{argmin}_i \ \ell_{\leq t-1,i}$ $\qquad\qquad$ ties broken uniformly

FPL($\eta$) $\qquad \widehat{i}_t = \text{argmin}_i \ \ell_{\leq t-1,i} + \frac{1}{\eta}\xi_{t,i}$ $\qquad$ indep. <u>additive</u> noise

Hedge($\eta$) $\quad w_i = \frac{e^{-\eta \ell_{\leq t-1,i}}}{Z}$ $\qquad\qquad$ Weighted Majority algorithm
for pred. with Expert Advice
Soft min

# Dropout

|            | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|------------|-------|-------|-------|-------|-------|
|            | 0     | ̸1    | 0     | 0     | ̸1    |
|            | 1     | 1     | 0     | 1     | 1     |
| day $t-1$  | 0     | 0     | ̸1    | ̸1    | 1     |

$\widehat{\ell}_{\leq t-1,i}$

# Dropout

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | $\not{1}$ | 0 | 0 | $\not{1}$ |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | $\not{1}$ | $\not{1}$ | 1 |
| $\widehat{\ell}_{\leq t-1,i}$ | 1 | 1 | 0 | 1 | 2 |

$$\widehat{\ell}_{t,i} = \beta_{t,i}\ell_{t,i}, \qquad \text{where } \beta_{t,i} \text{ iid Bernoilli}$$

# Dropout

|  | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ |
|---|---|---|---|---|---|
|  | 0 | $\cancel{1}$ | 0 | 0 | $\cancel{1}$ |
|  | 1 | 1 | 0 | 1 | 1 |
| day $t-1$ | 0 | 0 | $\cancel{1}$ | $\cancel{1}$ | 1 |
| $\widehat{\ell}_{\leq t-1,i}$ | 1 | 1 | 0 | 1 | 2 |

$\widehat{\ell}_{t,i} = \beta_{t,i}\ell_{t,i},$ where $\beta_{t,i}$ iid Bernoilli



FL on dropout
$$\widehat{i}_t = \operatorname*{argmin}_i \widehat{\ell}_{\leq t-1,i}$$

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

- FL is bad
- FPL($\eta$) and Hedge($\eta$) achieve optimal regret with tuning
  - fancy tunings: AdaHedge and Flipflop

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

- FL is bad
- FPL($\eta$) and Hedge($\eta$) achieve optimal regret with tuning
  - fancy tunings: AdaHedge and Flipflop
- FL on dropout requires no tuning

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

- FL is bad
- FPL($\eta$) and Hedge($\eta$) achieve optimal regret with tuning
  - fancy tunings: AdaHedge and Flipflop
- FL on dropout requires no tuning
  - dropout better noise for achieving optimal worst case regret
    - additive noise needs tuning - multiplicative noise does not
  - in iid case when gap between 1st and 2nd: $\log n$ regret

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

- FL is bad
- FPL($\eta$) and Hedge($\eta$) achieve optimal regret with tuning
  - fancy tunings: AdaHedge and Flipflop
- FL on dropout requires no tuning
  - dropout better noise for achieving optimal worst case regret
    - additive noise needs tuning - multiplicative noise does not
  - in iid case when gap between 1st and 2nd: $\log n$ regret
- In the meantime
  - new fancy algorithms by
  Haipeng Luo, Rob Schapire & Tim van Erven, Wouter Koolen

# How good?

Optimal worst case regret: $\sqrt{L^* \ln n} + \ln n$

- FL is bad
- FPL($\eta$) and Hedge($\eta$) achieve optimal regret with tuning
  - fancy tunings: AdaHedge and Flipflop
- FL on dropout requires no tuning
  - dropout better noise for achieving optimal worst case regret
    - additive noise needs tuning - multiplicative noise does not
  - in iid case when gap between 1st and 2nd: $\log n$ regret
- In the meantime
  - new fancy algorithms by
  Haipeng Luo, Rob Schapire & Tim van Erven, Wouter Koolen
  - also no tuning, many other advantages

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t \quad \longrightarrow \quad$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t \quad \longrightarrow \quad$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta\ell_{\leq t-1,i}}}{Z} \quad \longrightarrow \quad$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta\boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t$ $\longrightarrow$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t$ $\longrightarrow$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta\ell_{\leq t-1,i}}}{Z}$ $\longrightarrow$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta\boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

- FL minimum eigenvector calculation of $\boldsymbol{L}_{\leq t-1}$: $O(n^2)$

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t \quad \longrightarrow \quad$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t \quad \longrightarrow \quad$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta \ell_{\leq t-1,i}}}{Z} \quad \longrightarrow \quad$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta \boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

- FL minimum eigenvector calculation of $\boldsymbol{L}_{\leq t-1}$: $\quad O(n^2)$
- Is there $O(n^2)$ perturbation with optimal regret bound?

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t \quad \longrightarrow \quad$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t \quad \longrightarrow \quad$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta \ell_{\leq t-1,i}}}{Z} \quad \longrightarrow \quad$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta \boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

- FL minimum eigenvector calculation of $\boldsymbol{L}_{\leq t-1}$: $\quad O(n^2)$
- Is there $O(n^2)$ perturbation with optimal regret bound?

- Follow the skipping leader:
  Drop entire loss $\boldsymbol{L}_t$ with probability $\frac{1}{2}$
  = Online Bagging

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t \quad \longrightarrow \quad$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t \quad \longrightarrow \quad$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta \ell_{\leq t-1,i}}}{Z} \quad \longrightarrow \quad$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta \boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

- FL minimum eigenvector calculation of $\boldsymbol{L}_{\leq t-1}$: $O(n^2)$
- Is there $O(n^2)$ perturbation with optimal regret bound?

- Follow the skipping leader:
  Drop entire loss $\boldsymbol{L}_t$ with probability $\frac{1}{2}$
  = Online Bagging
- Proof techniques break down
  - settled for vector case and independent multiplicative noise
  = dropout

# Our path to dropout

- Loss vectors $\boldsymbol{\ell}_t \longrightarrow$ loss matrices $\boldsymbol{L}_t$
- Prob. vectors $\boldsymbol{w}_t \longrightarrow$ density matrices $\mathbf{W}_t$
- Hedge $w_{t,i} = \frac{e^{-\eta \ell_{\leq t-1,i}}}{Z} \longrightarrow$ Matrix Hedge

$$\mathbf{W}_t = \frac{\exp\left(-\eta \boldsymbol{L}_{\leq t-1}\right)}{Z'}$$

- Matrix Hedge $O(n^3)$ per update

- FL minimum eigenvector calculation of $\boldsymbol{L}_{\leq t-1}$: $O(n^2)$
- Is there $O(n^2)$ perturbation with optimal regret bound?

- Follow the skipping leader:
  Drop entire loss $\boldsymbol{L}_t$ with probability $\frac{1}{2}$
  = Online Bagging
- Proof techniques break down
  - settled for vector case and independent multiplicative noise
  = dropout
- Follow the skipping leader has linear regret   [Lugosi,Neu2014]

# What regularization?

Hedge($\eta$)     relative entropy

# What regularization?

| Hedge($\eta$) | relative entropy |
|---|---|
| FPL($\eta$) | additive $\frac{1}{\eta}$ log exponential noise $=$ Hedge($\eta$) |

# What regularization?

Hedge($\eta$)        relative entropy
FPL($\eta$)          additive $\frac{1}{\eta}$ log exponential noise $=$ Hedge($\eta$)

FL on dropout   tricky

Feed forward NN                     [Wagner, Wang, Liang 2013]
Logistic regression                    [Helmbold, Long 2014]
Linear loss case                          [ALST 2014]

# Outline

## Simple algorithms

Any deterministic alg. (such as FL) has huge regret

- For $T$ trials: give algorithm's expert a unit of loss
- Loss of alg.: $T$      loss of best: $\leq \frac{T}{n}$

# Simple algorithms

Any deterministic alg. (such as FL) has huge regret

- For $T$ trials: give algorithm's expert a unit of loss
- Loss of alg.: $T$     loss of best: $\leq \frac{T}{n}$

regret: $\geq \underbrace{T}_{nL^*} - \underbrace{\frac{T}{n}}_{L^*} = (n-1)L^*$

# Simple algorithms

Any deterministic alg. (such as FL) has huge regret

- For $T$ trials: give algorithm's expert a unit of loss
- Loss of alg.: $T$     loss of best: $\leq \frac{T}{n}$

regret: $\geq \underbrace{T}_{nL^*} - \underbrace{\frac{T}{n}}_{L^*} = (n-1)L^*$

Recall optimum regret: $\sqrt{L^* \ln n} + \ln n$

FL with random ties

# Simple algorithms

Any deterministic alg. (such as FL) has huge regret

- For $T$ trials: give algorithm's expert a unit of loss
- Loss of alg.: $T$        loss of best: $\leq \frac{T}{n}$

  regret: $\geq \underbrace{T}_{nL^*} - \underbrace{\frac{T}{n}}_{L^*} = (n-1)L^*$

Recall optimum regret: $\sqrt{L^* \ln n} + \ln n$

FL with random ties

- Give every expert one unit of loss
  - iterate $L^* + 1$ times
- Loss per sweep        $\frac{1}{n} + \frac{1}{n-1} + \ldots + \frac{1}{2} + 1 \approx \ln n$
- Loss of alg.: $(L^* + 1) \ln n$    loss of best: $L^*$
  regret: $L^* \ln n$

**Unit rule**

- Adversary forces more regret by splitting loss vectors into units

$$\begin{pmatrix} \textcolor{red}{1} \\ 0 \\ \textcolor{blue}{1} \\ \textcolor{green}{1} \end{pmatrix} \longrightarrow \begin{pmatrix} \textcolor{red}{1} \\ 0 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \textcolor{blue}{1} \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ \textcolor{green}{1} \end{pmatrix}$$

**Unit rule**

- Adversary forces more regret by splitting loss vectors into units

$$\begin{pmatrix} \mathbf{1} \\ 0 \\ \mathbf{1} \\ \mathbf{1} \end{pmatrix} \longrightarrow \begin{pmatrix} \mathbf{1} \\ 0 \\ 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \mathbf{1} \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ \mathbf{1} \end{pmatrix}$$

**Swapping rule**                                   $\ell_{\leq T,i}$

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $E_1$ | 1 | 1 | 1 | 1 | 1 | 1 | **1** | 1 | 1 | | 9 |
| $E_2$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | | 8 |
| $E_3$ | 1 | 1 | 1 | **1** | 1 | 1 | 1 | 1 | 1 | 1 | 10 |
| $E_4$ | 1 | 1 | 1 | 1 | 1 | 1 | | | | | 6 |

**Unit rule**

- Adversary forces more regret by splitting loss vectors into units

$$\begin{pmatrix}\mathbf{1}\\0\\\mathbf{1}\\\mathbf{1}\end{pmatrix} \longrightarrow \begin{pmatrix}\mathbf{1}\\0\\0\\0\end{pmatrix} \begin{pmatrix}0\\0\\\mathbf{1}\\0\end{pmatrix} \begin{pmatrix}0\\0\\0\\\mathbf{1}\end{pmatrix}$$

**Swapping rule**                                         $\ell_{\leq T,i}$

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $E_1$ | 1 | 1 | 1 | 1 | 1 | 1 | **1** | 1 | 1 | | 9 |
| $E_2$ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | | | 8 |
| $E_3$ | 1 | 1 | 1 | **1** | 1 | 1 | 1 | 1 | 1 | 1 | 10 |
| $E_4$ | 1 | 1 | 1 | 1 | 1 | 1 | | | | | 6 |

- 1's occur in some order
- Worst case: **1** before **1**
- Otherwise adversary benefits from swapping

# Worst-case pattern

```
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     1     1     1     1     1     1     1
 1     1     1     1     1     1     1     1     1
```

# Cost per sweep

Assume we have $s$ leaders

# Cost per sweep

Assume we have $s$ leaders

$s$ leader get unit
ignore non-leaders
$$\left\{ \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} \right.$$

Assume we have $s$ leaders

$s$ leader get unit
ignore non-leaders
$\left\{ \begin{array}{l} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} \right.$

**FL**

$$\frac{1}{s} + \frac{1}{s-1} + \frac{1}{s-2} + \frac{1}{s-3} + \ldots + \underbrace{\frac{1}{s - s - 2}}_{2} + \underbrace{\frac{1}{s - s - 1}}_{1}$$

$\approx \ln s$

# Cost per sweep

Assume we have $s$ leaders

$s$ leader get unit
ignore non-leaders
$\left\{\begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array}\right.$

**FL**

$$\frac{1}{s} + \frac{1}{s-1} + \frac{1}{s-2} + \frac{1}{s-3} + \ldots + \underbrace{\frac{1}{s-s-2}}_{2} + \underbrace{\frac{1}{s-s-1}}_{1}$$

$\approx \ln s$

**Dropout**

$$\frac{1}{s} + \frac{1}{s-1/2} + \frac{1}{s-2/2} + \frac{1}{s-3/2} + \ldots + \frac{1}{s-(s-2)/2} + \frac{1}{s-(s-1)/2}$$

Assume we have $s$ leaders

$s$ leader get unit
ignore non-leaders
$$\left\{ \begin{array}{l} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} \right.$$

**FL**

$$\frac{1}{s} + \frac{1}{s-1} + \frac{1}{s-2} + \frac{1}{s-3} + \ldots + \underbrace{\frac{1}{s-s-2}}_{2} + \underbrace{\frac{1}{s-s-1}}_{1}$$

$$\approx \ln s$$

**Dropout**

$$\frac{2}{2s} + \frac{2}{2s-1} + \frac{2}{2s-2} + \frac{2}{s-3} + \ldots + \frac{2}{2s-(s-2)} + \frac{2}{2s-(s-1)}$$

$$\approx 2\left(\ln 2s - \ln s\right) = 2\ln 2$$

# Cost per sweep

Assume we have $s$ leaders

$s$ leader get unit
ignore non-leaders
$$\left\{\begin{array}{l} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array}\right.$$

**FL**

$$\frac{1}{s} + \frac{1}{s-1} + \frac{1}{s-2} + \frac{1}{s-3} + \ldots + \underbrace{\frac{1}{s-s-2}}_{2} + \underbrace{\frac{1}{s-s-1}}_{1}$$

$$\approx \ln s$$

**Dropout**

$$\frac{2}{2s} + \frac{2}{2s-1} + \frac{2}{2s-2} + \frac{2}{s-3} + \ldots + \frac{2}{2s-(s-2)} + \frac{2}{2s-(s-1)}$$

$$\approx 2\left(\ln 2s - \ln s\right) = 2 \ln 2$$

FL

- One sweep

$$\frac{1}{n} + \frac{1}{n-1} + \ldots + \frac{1}{2} + 1 \approx (\ln n) - 1$$

- Optimal

# $L^* = 0$ - one expert incurs no loss

FL

- One sweep

$$\frac{1}{n} + \frac{1}{n-1} + \ldots + \frac{1}{2} + 1 \approx (\ln n) - 1$$

- Optimal

Dropout

- \# of leaders reduced by half in each sweep
- $\approx \log_2 n$ sweeps    times    $\leq 2 \ln 2 = 1.386$
  ========================
  $2 \ln n$

- Focus on first $L$ sweeps
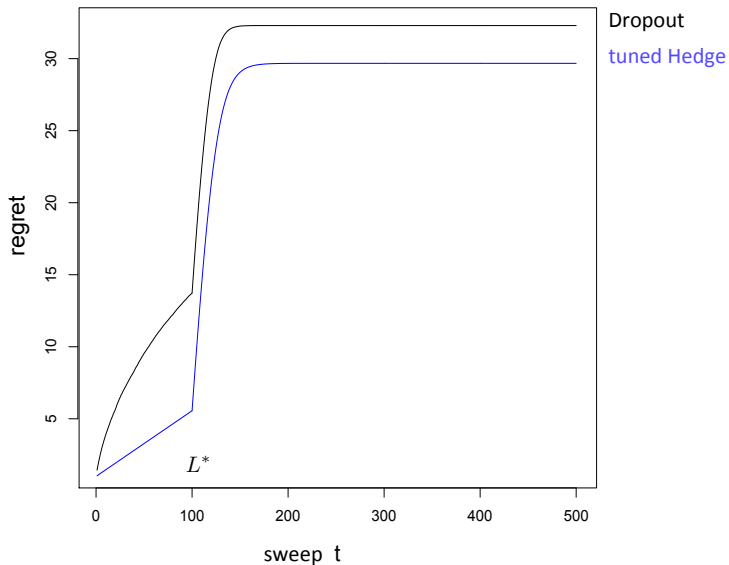- Only occurs constant regret if number of leaders $> 1$

## Overview of proof for noisy case

- Focus on first $L$ sweeps
- Only occurs constant regret if number of leaders $> 1$

- Prob. that number of leaders $> 1$ is at most $\sqrt{\frac{\ln n}{q+1}}$ for sweep $q$

# Overview of proof for noisy case

- Focus on first $L$ sweeps
- Only occurs constant regret if number of leaders $> 1$

- Prob. that number of leaders $> 1$ is at most $\sqrt{\frac{\ln n}{q+1}}$ for sweep $q$

- For Hedge$(\eta)$ and FPL$(\eta)$ cost per sweep constant and dependent on $\eta$

# Outlook

- Combinatorial experts
- Matrix case
- Where else can dropout perturbations be used?
- Dropout for convex losses
- Dropout for neural nets

- Combinatorial experts
- Matrix case
- Where else can dropout perturbations be used?
- Dropout for convex losses
- Dropout for neural nets
- Privacy

# [Lugosi, Neu 2014] dense counter example

| 0 | 1* |
|---|-----|
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| $\frac{n-1}{n}$ | $\frac{1}{2}$ |

Iterate this pattern $n$ times:

$$\sum_{i=1}^{n} \left( \frac{n-i}{n-i+1} + \frac{1}{2} \right)$$

$$\approx n - \ln n + \frac{n}{2}$$

$L^* = n$: Follow the Scipping Leader has linear regret

# How does dropout ovoid this example?

| | |
|---|---|
| 0 | 1 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| $\frac{n-1}{n}$ | $\frac{1}{\frac{n-1}{2}}$ |

# How does dropout ovoid this example?

| | |
|---|---|
| 0 | 1 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| 1 | 0 |
| $\frac{n-1}{n}$ | $\frac{1}{\frac{n-1}{2}}$ |

It leaves the adversary clueless as to who the leader is
i.e. privacy against adversary

## sparse counter example

| | |
|---|---|
| $0$ | $1^*$ |
| $\frac{1}{n-1}$ | $0$ |
| $\frac{1}{n-1}$ | $0$ |
| $\frac{1}{n-1}$ | $0$ |
| $\frac{1}{n-1}$ | $0$ |
| $\frac{1}{n-1}$ | $0$ |
| $\frac{n-1}{n^2}$ | $\frac{1}{2}$ |

Iterate this pattern $n$ times:

$$\sum_{i=1}^{n} \left( \frac{n-i}{(n-i+1)^2} + \frac{1}{2} \right)$$
$$= \sum_{i=1}^{n} \left( \frac{1}{n-i+1} - \frac{1}{(n-i+1)^2} + \frac{1}{2} \right)$$
$$\approx \ln n - O(1) + \frac{n}{2}$$

$L^* = \ln n$: Follow the Scipping Leader has linear regret