

# Fundamentals

Seminar Data Analytics I  
International Master's Program in Data Analytics  
University of Hildesheim  
Summer Semester 2018

Daniel Obando Montero  
Shabanaz Chamurally  
Daniela Thyssens

# Motivation

## 1. Stochastic Gradient Descent Training for L1-regularized Log-linear Models with Cumulative Penalty

- Stochastic gradient descent (SGD) is attractive because it often requires much less training time in practice than batch training algorithms.
- However, L1-regularization cannot be efficiently applied in SGD training.
- Paper present a simple method for solving two problems in SGD learning.

## 2. Curriculum Learning

- Humans and animals learn much better when the examples are presented in order of increasing difficulty and organized in a meaningful order.
- The paper formalizes how a curriculum can be used in the context of machine learning.

## 3. Combined Regression & Ranking

- Many real-world data mining tasks require the achievement of two distinct goals when applied to unseen data: first, to induce an accurate preference ranking, and second to give good regression performance.
- The paper present an efficient and effective Combined Regression and Ranking method (CRR) that optimizes regression and ranking objectives simultaneously.

# Stochastic Gradient Descent Training for L1-regularized Log-linear Models with Cumulative Penalty

Yoshimasa Tsuruoka, Jun'ichiTsuji and Sophia Ananiadou

ACL '09 Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language  
Processing of the AFNLP: Volume 1 - Volume 1 Pages 477-485  
Suntec, Singapore — August 02 - 07, 2009

Seminar Data Analytics I  
International Master's Program in Data Analytics  
University of Hildesheim  
Summer Semester 2018

Daniel Obando Montero

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. Experiments
6. Criticism
7. Conclusions

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. Experiments
6. Criticism
7. Conclusions

# Introduction

- SGD is a very attractive learning framework.
- Not efficient with L1.
- 2 problems:
  - Inefficiency of applying the L1 penalty to the weights of the features.
  - Doesn't always lead to compact models.

# Introduction

- Objective: present a simple method for solving the two problems in SGD learning.

# Introduction

- How?
  - keep track of the total penalty and the penalty applied on each weight
  - The penalty is applied based on the difference between those values.

# Introduction

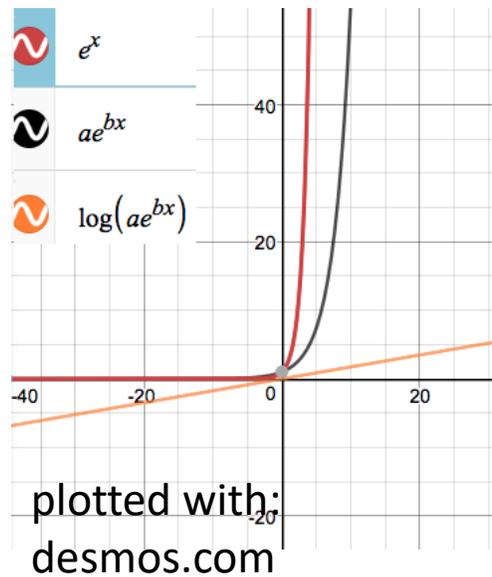
- How?
  - vs “state-of-the-art” OWL-QN
  - experiments with:
    - text chunking, named entity recognition and POS tagging

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. Experiments
6. Criticism
7. Conclusions

# Log-Linear Models in NLP and L1

- Where  $f_i$  indicates the occurrence of feature  $i$  and  $w_i$  is the weight of the feature.
- $Z(\mathbf{x})$  is a partition function that allows a well formed normal distribution.



$$p(\mathbf{y}|\mathbf{x}) = \frac{e^{\sum_i w_i f_i(\mathbf{y}, \mathbf{x})}}{Z(\mathbf{x})}$$

$$Z(\mathbf{x}) = \sum_y e^{\sum_i w_i f_i(\mathbf{y}, \mathbf{x})}$$

(Tsuruoka,  
2009)

# Log-Linear Models in NLP and L1

- Maximize the conditional likelihood of the data.
- $N$  = # of training samples.
- $y_j$  is the correct output of input  $x_j$
- $R$  is the L1-regularization term known as LASSO regression (Least Absolute Shrinkage and Selection Operator).
- $C$  is the metaparameter that controls the degree of regularization.

$$\mathcal{L}_{\mathbf{w}} = \sum_{j=1}^N \log p(\mathbf{y}_j | \mathbf{x}_j; \mathbf{w}) - R(\mathbf{w})$$

$$R(\mathbf{w}) = C \sum_i |w_i|$$

$$L(j, \mathbf{w}) \forall \text{ sample } \in \log p(\mathbf{y}_j | \mathbf{x}_j; \mathbf{w})$$

$$\mathcal{L}_{\mathbf{w}} = \sum_{j=1}^N L(j, \mathbf{w}) - C \sum_i |w_i|$$

(Tsuruoka, 2009)

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. **SGD - L1**
4. OWL-QN
5. Experiments
6. Criticism
7. Conclusions

# Stochastic Gradient Descent(SGD)

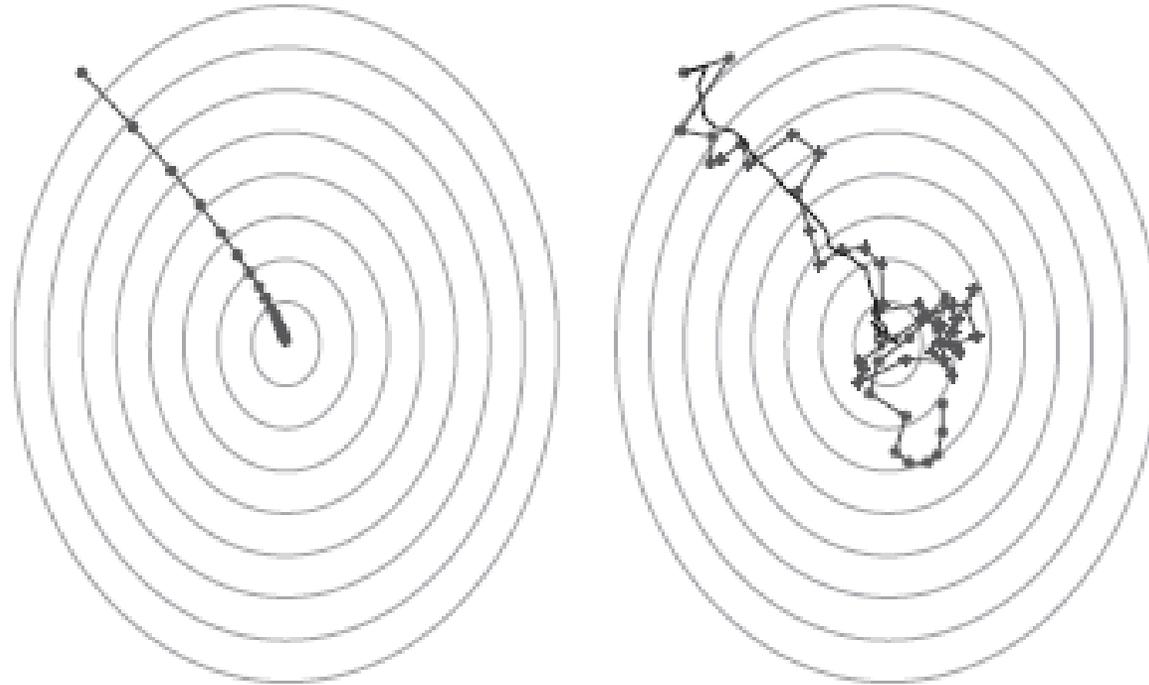
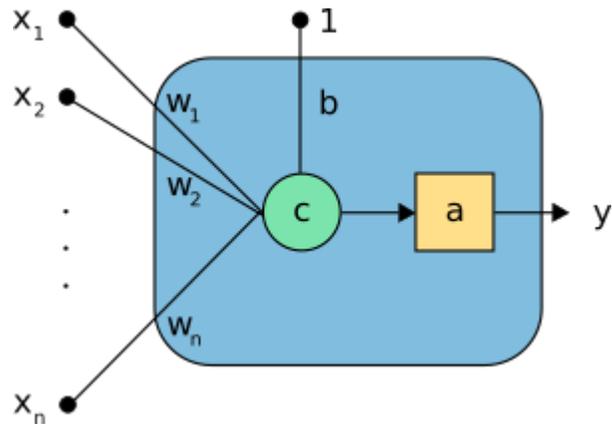


figure source: researchgate.net

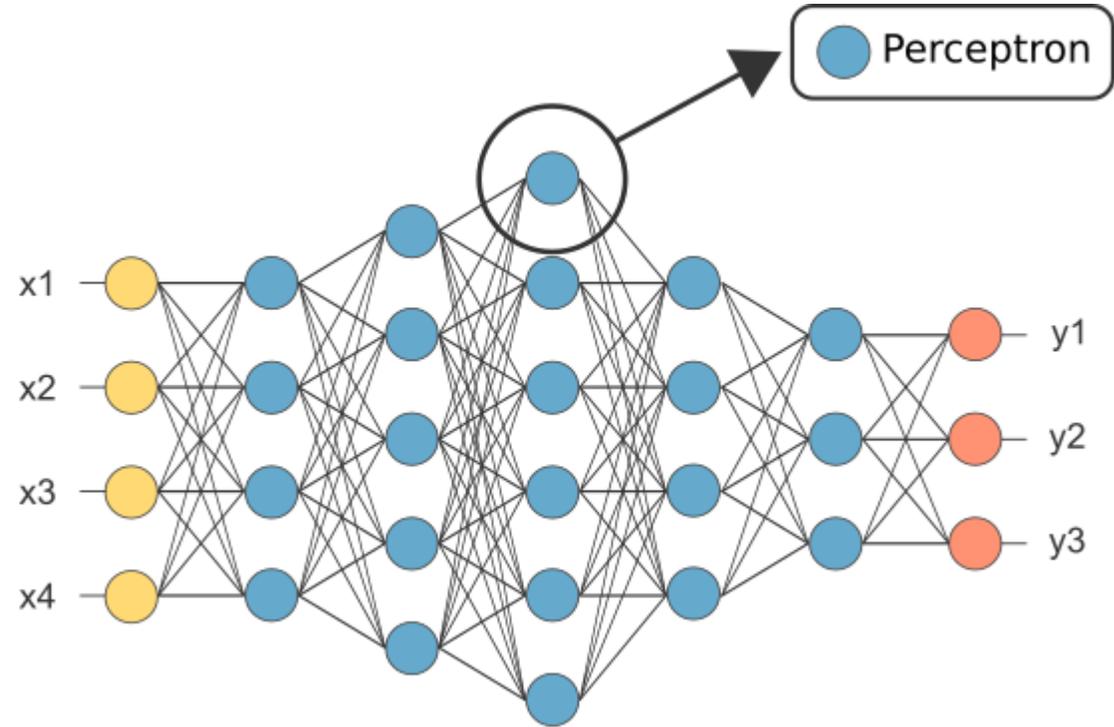
# Perceptron



$$h(x) = Wx + b$$

$$y = a(h)$$

figures source: neuraldesigner.com

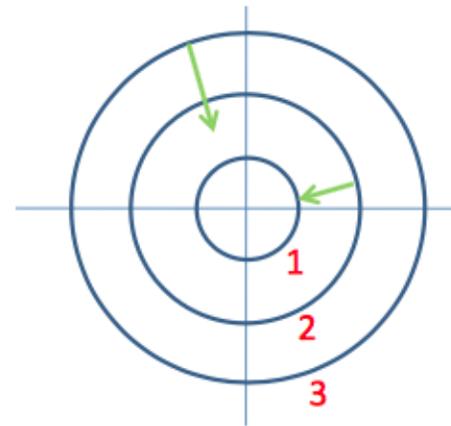


# SGD - L1

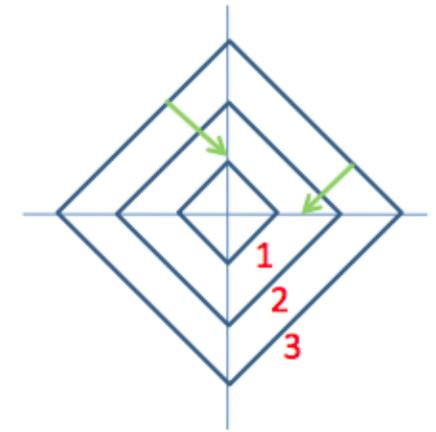
- SGD-L1 (Naive).
- SGD-L1 (Clipping + Lazy Update).
- SGD-L1 (Cumulative).
- SGD-L1 (Cumulative + Exponential Decay).

# SGD - L1

- Why L1?
- Is capable of learning good models when most features are irrelevant.
- Many of the parameters are exactly zero.
- L2 “pushes” to 0 but it never reaches it



Negative gradient of  $L_2$ -norm  
always points directly toward 0



“Negative gradient” of  $L_1$ -norm  
(direction of steepest descent)  
points toward coordinate axes

(Andrew, 2007)

# SGD and L1

The updates of the weights of the features is given by the following formula at the given  $j$  sample.

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \eta k \frac{\partial}{\partial \mathbf{w}} (L(j, \mathbf{w}) - \frac{C}{N} \sum_i |w_i|)$$

Then, the update for the weight of each feature  $i$  is:

$$w_i^{k+1} = w_i^k + \eta k \frac{\partial}{\partial w_i} (L(j, \mathbf{w}) - \frac{C}{N} |w_i|)$$

(Tsuruoka, 2009)

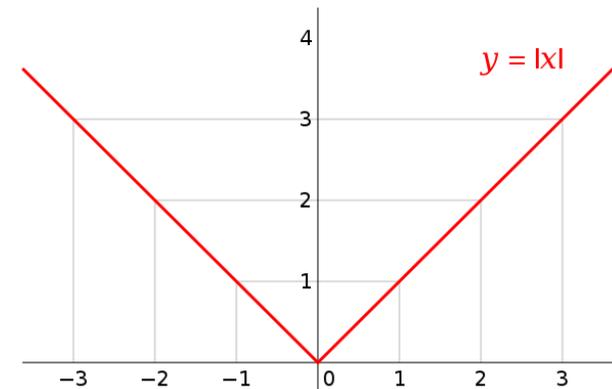


figure source: wikipedia.org

# SGD and L1 (naive)

$$w_i^{k+1} = w_i^k + \eta^k \frac{\partial L(j, \mathbf{w})}{\partial w_i} - \frac{C}{N} \eta_k \text{sign}(w_i)$$

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & x < 0 \\ 0 & x = 0 \end{cases}$$

(Tsuruoka, 2009)

# SGD and L1 (naive)

- Two problems:
  - at each update, the penalty must be applied to all features.
  - it does not produce a compact model.
    - i.e. most of the weights of the features do not become zero as a result of training.

# SGD-L1 (Clipping + Lazy Update)

$$w_i^{k+\frac{1}{2}} = w_i^k + \eta_k \frac{\partial L(j, \mathbf{w})}{\partial w_i} \Bigg|_{w=w^k}$$

**if**  $w_i^{k+\frac{1}{2}} > 0$  **then**

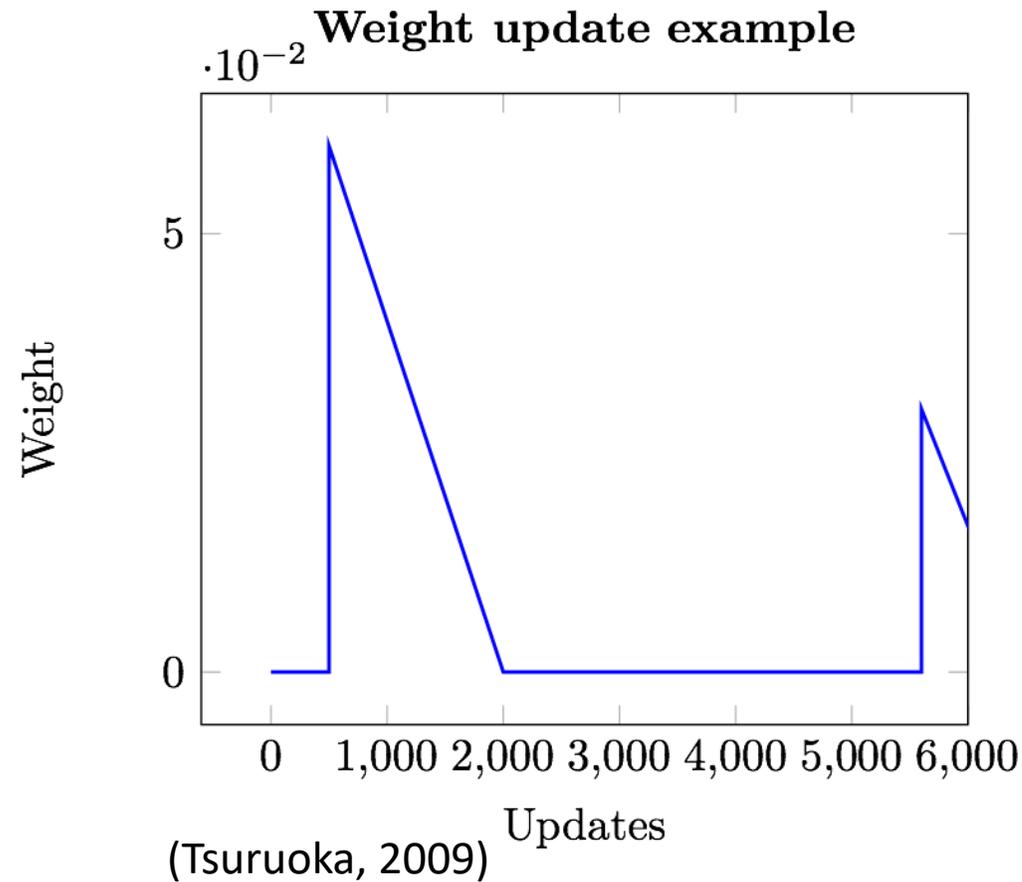
$$\left| w_i^{k+1} = \max(0, w_i^{k+\frac{1}{2}} - \frac{C}{N} \eta_k) \right.$$

**else if**  $w_i^{k+\frac{1}{2}} < 0$  **then**

$$\left| w_i^{k+1} = \min(0, w_i^{k+\frac{1}{2}} + \frac{C}{N} \eta_k) \right.$$

(Carpenter, 2008)

# SGD-L1 (Clipping + Lazy Update)



# SGD-L1 (Cumulative)

$$u_k = \frac{C}{N} \sum_{t=1}^k \eta_t \quad q_i^k = \sum_{t=1}^k (w_i^{t+1} - w_i^{t+\frac{1}{2}})$$

- $u_k$  is the value of the total penalty each weight could have received.
- $q_i^k$  is the total L1 penalty that  $w_i$  has actually received.

$$w_i^{k+\frac{1}{2}} = w_i^k + \eta_k \left. \frac{\partial L(j, \mathbf{w})}{\partial w_i} \right|_{w=w^k}$$

**if**  $w_i^{k+\frac{1}{2}} > 0$  **then**

$$\left| w_i^{k+1} = \max(0, w_i^{k+\frac{1}{2}} - (u_k + q_i^{k-1})) \right.$$

**else if**  $w_i^{k+\frac{1}{2}} < 0$  **then**

$$\left| w_i^{k+1} = \min(0, w_i^{k+\frac{1}{2}} + (u_k - q_i^{k-1})) \right.$$

(Tsuruoka, 2009)

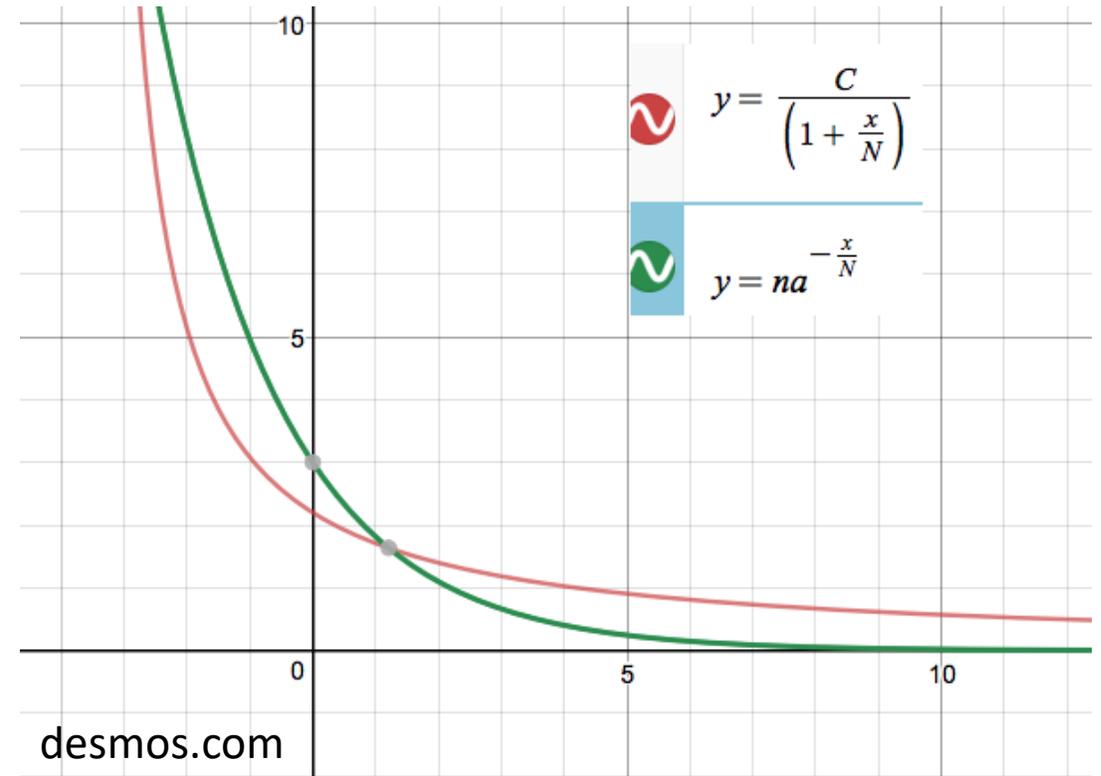
# SGD-L1 (Cumulative + Exponential Decay)

$\eta_0$  and alpha are meta-parameters

$$\eta_k = \frac{\eta_0}{1 + \frac{k}{N}}$$

$$\eta_k = \eta_0 \alpha^{-\frac{k}{N}}$$

(Tsuruoka, 2009)



```

1 train ( $C$ ):
2    $u \leftarrow 0$ 
3   Initialize  $w_i$  and  $q_i$  with zero for all  $i$ 
4   for  $k = 0$  to  $maxIterations$  do
5      $\eta \leftarrow \text{learningRate}(k)$ 
6      $u \leftarrow u + \eta \frac{C}{N}$ 
7     Select sample  $j$  randomly
8     updateWeights( $j$ )
9   end
10
11 updateWeights ( $j$ ):
12   for  $i \in \text{features used in sample } j$  do
13      $w_i \leftarrow w_i + \eta \frac{\partial L(j, \mathbf{w})}{\partial w_i}$ 
14     applyPenalty( $i$ )
15   end
16
17 applyPenalty ( $i$ ):
18    $z \leftarrow w_i$ 
19   if  $w_i > 0$  then
20      $w_i = \max(0, w_i - (u + q_i))$ 
21   else if  $w_i < 0$  then
22      $w_i = \min(0, w_i + (u - q_i))$ 
23    $q_i \leftarrow q_i + (w_i - z)$ 

```

(Tsuruoka, 2009)

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. **OWL-QN**
5. Experiments
6. Criticism
7. Conclusions

# OWL-QN

- Orthant-Wise Limited-memory Quasi Newton
- Scalable training of L1-regularized log-linear models
- Galen Andrew and Jianfeng Gao. 2007.
- Proposal:
  - An algorithm based on L-BFGS for training large-scale log-linear models using L1-regularization.

# OWL-QN

Newton

Gradient Descent

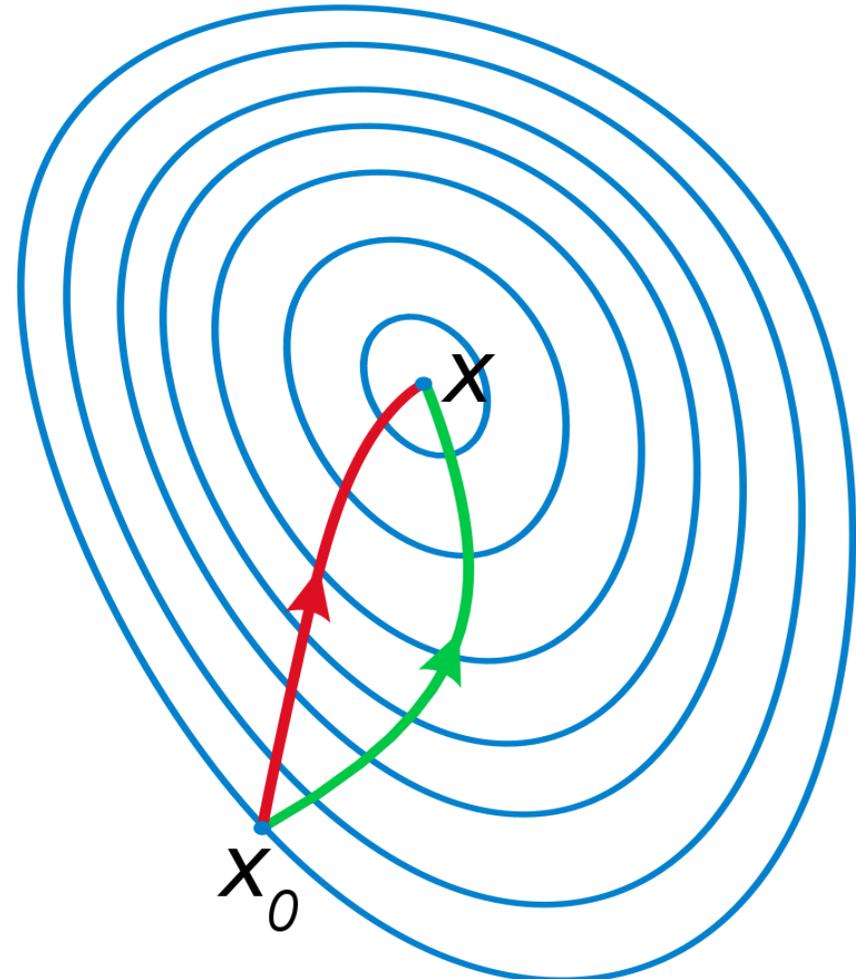


figure source: wikipedia.org

# OWL-QN

- Newton's method: Jacobian/Hessian.
- Quasi Newton uses approximations.
- Broyden-Fletcher-Goldfarb-Shanno (BFGS).
- Convergence and quadratic Taylor Series.
- Limited-memory BFGS (L-BFGS).

# OWL-QN

- $\mathbf{H}_k$ : L-BFGS approximation to the inverse Hessian of the loss.

$$f(x) = \ell(x) + r(x)$$

$$r(x) = C \sum_i |x_i|, \quad C > 0$$

$$\ell(x) = - \sum_{j=1}^M \log P_w(y_j | x_j)$$

$$x^{k+1} = x^k - \alpha \mathbf{H}_k g^k, \quad \text{for } \alpha \in (0, \infty)$$

(Andrew, 2007)

# OWL-QN

```
1 choose initial point  $x^0$ 
2  $displacement \leftarrow \{\}$ ,  $gradientChange \leftarrow \{\}$ 
3 for  $k = 0$  to  $maxIterations$  do
4   Compute  $orthant^k = -\diamond f(x^k)$ 
5   Compute  $descentDirection^k \leftarrow$ 
      $\mathbf{H}_k orthant^k$  using  $displacement$  and  $gradientChange$ 
6    $projection^k \leftarrow \pi(descentDirection^k; orthant^k)$ 
7   Find  $x^{k+1}$  with constrained line search
8   if termination condition satisfied then
9     | Stop and return  $x^{k+1}$ 
10  Update  $displacement$  with  $s^k = x^{k+1} - x^k$ 
11  Update  $gradientChange$  with  $y^k = \nabla \ell(x^{k+1}) - \nabla \ell(x^k)$ 
12 end
```

(Andrew, 2007)

- Variant of L-BFGS
- pseudo gradients.
- projection of descent onto the defined orthant.
- to remain in the valid region were the projection is valid.
- discarding the vectors from last iteration.

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. **Experiments**
6. Criticism
7. Conclusions

# Experiments

- Conditional Random Fields (CRF):
  - Text Chunking
  - Named Entity Recognition
  - Part-Of-Speech Tagging (POS-Tagging)

# Experiments

- Processor Xeon 2.13GHz
- CRF++ version 0.50 developed by Taku Kudo

# Experiments - Text Chunking

- Training data: 8936 sentences.
- Held-out data: 1000
- C was tuned to maximize the likelihood of the held-out data.
- The learning rate parameters were tuned to maximize the value of the objective function in 30 passes.

# Experiments - Text Chunking

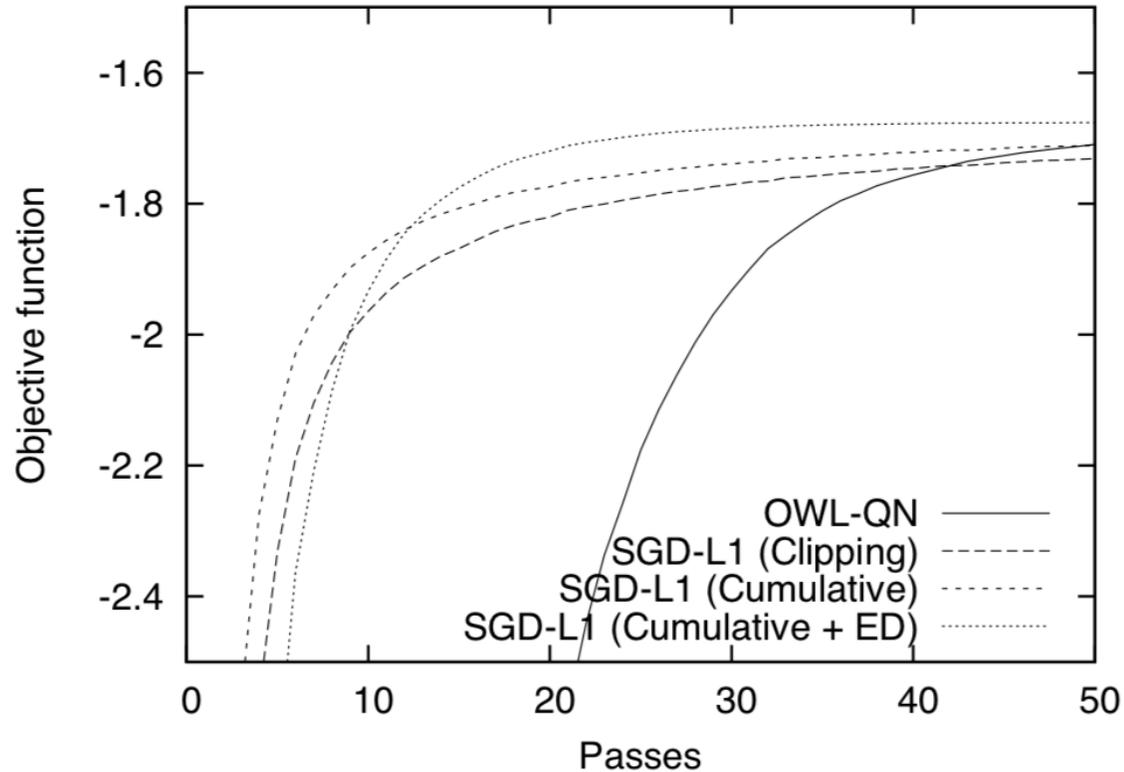


Figure 3: CoNLL 2000 chunking task: Objective (Tsuruoka, 2009)

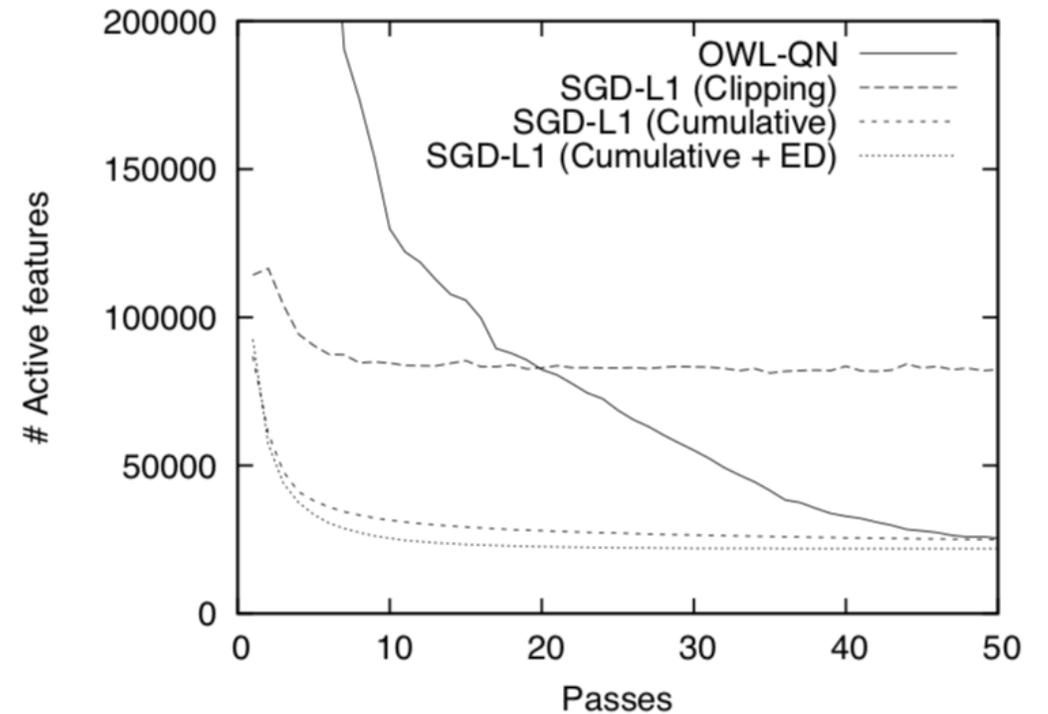


Figure 4: CoNLL 2000 chunking task: Number of active features.

# Experiments - Text Chunking

	Passes	$\mathcal{L}_w/N$	# Features	Time (sec)	F-score
OWL-QN	160	-1.583	18,109	598	93.62
SGD-L1 (Naive)	30	-1.671	455,651	1,117	93.64
SGD-L1 (Clipping + Lazy-Update)	30	-1.671	87,792	144	93.65
SGD-L1 (Cumulative)	30	-1.653	28,189	149	93.68
SGD-L1 (Cumulative + Exponential-Decay)	30	-1.622	23,584	148	93.66

Table 1: CoNLL-2000 Chunking task. Training time and accuracy of the trained model on the test data.

(Tsuruoka, 2009)

# Experiments - Named Entity Recognition

- 18,546 sentences

	Passes	$\mathcal{L}_w/N$	# Features	Time (sec)	F-score
OWL-QN	161	-2.448	30,710	2,253	71.76
SGD-L1 (Naive)	30	-2.537	1,032,962	4,528	71.20
SGD-L1 (Clipping + Lazy-Update)	30	-2.538	279,886	585	71.20
SGD-L1 (Cumulative)	30	-2.479	31,986	631	71.40
SGD-L1 (Cumulative + Exponential-Decay)	30	-2.443	25,965	631	71.63

Table 2: NLPBA 2004 Named entity recognition task. Training time and accuracy of the trained model on the test data.

(Tsuruoka, 2009)

# Experiments - Named Entity Recognition

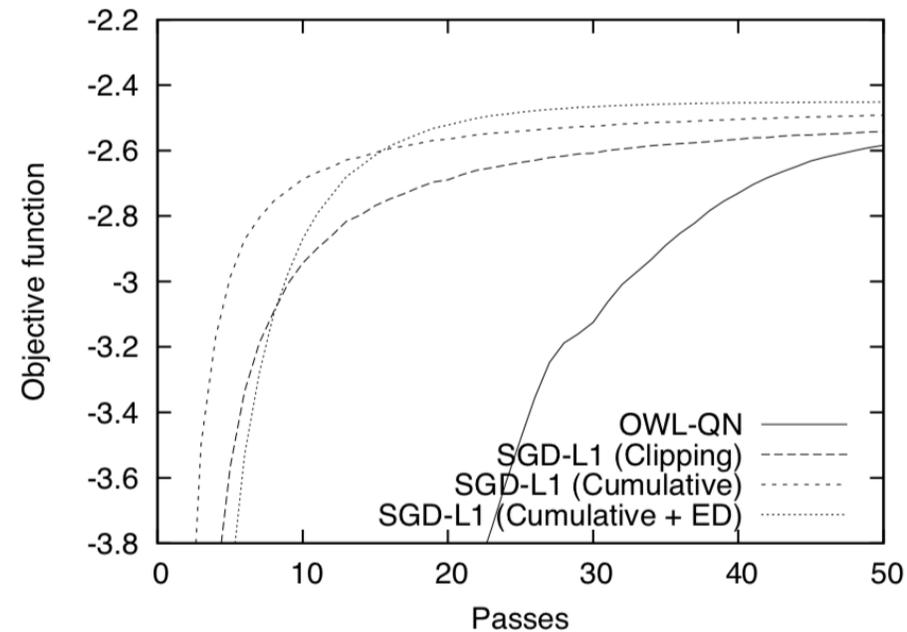


Figure 5: NLPBA 2004 named entity recognition task: Objective.

(Tsuruoka, 2009)

# Experiments - POS-Tagging

	Passes	$\mathcal{L}_w/N$	# Features	Time (sec)	Accuracy
OWL-QN	124	-1.941	50,870	5,623	97.16%
SGD-L1 (Naive)	30	-2.013	2,142,130	18,471	97.18%
SGD-L1 (Clipping + Lazy-Update)	30	-2.013	323,199	1,680	97.18%
SGD-L1 (Cumulative)	30	-1.987	62,043	1,777	97.19%
SGD-L1 (Cumulative + Exponential-Decay)	30	-1.954	51,857	1,774	97.17%

Table 3: POS tagging on the WSJ corpus. Training time and accuracy of the trained model on the test data.

(Tsuruoka, 2009)

# Experiments - POS-Tagging

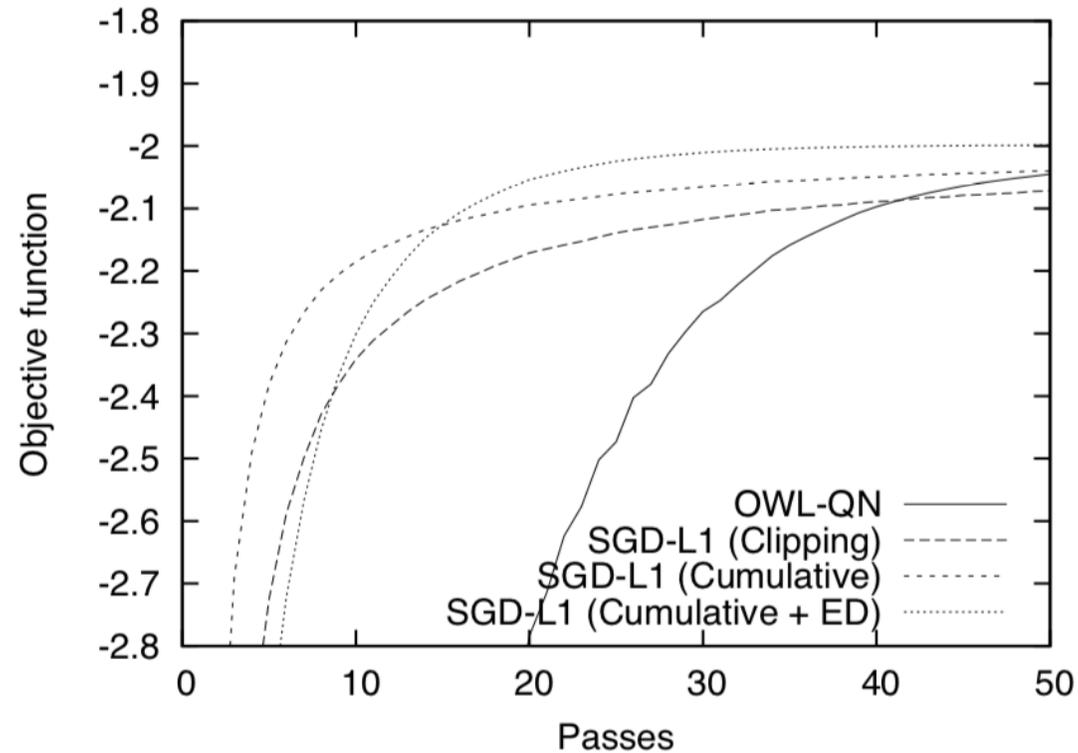


Figure 6: POS tagging task: Objective.

(Tsuruoka, 2009)

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. Experiments
6. **Criticism**
7. Conclusions

# Criticism

- Google Scholar references since 2009: 161
- “State-of-the-art” Google Scholar references since 2007: 459

## Stochastic gradient descent training for l1-regularized log-linear models with cumulative penalty

[Y Tsuruoka](#), [J Tsujii](#), [S Ananiadou](#) - ... of the Joint Conference of the 47th ..., 2009 - dl.acm.org

Stochastic gradient descent (SGD) uses approximate gradients estimated from subsets of the training data and updates the parameters in an online fashion. This learning framework is attractive because it often requires much less training time in practice than batch training algorithms. However, L1-regularization, which is becoming popular in natural language processing because of its ability to produce compact models, cannot be efficiently applied in SGD training, due to the large dimensions of feature vectors and the fluctuations of ...

☆ [Cited by 164](#) [Related articles](#) [All 19 versions](#)

## Scalable training of L 1-regularized log-linear models

[G Andrew](#), [J Gao](#) - Proceedings of the 24th international conference on ..., 2007 - dl.acm.org

The L-BFGS limited-memory quasi-Newton method is the algorithm of choice for optimizing the parameters of large-scale log-linear models with L2 regularization, but it cannot be used for an L1-regularized loss due to its non-differentiability whenever some parameter is zero. Efficient algorithms have been proposed for this task, but they are impractical when the number of parameters is very large. We present an algorithm Orthant-Wise Limited-memory Quasi-Newton (OWL-QN), based on L-BFGS, that can efficiently optimize the L1-regularized ...

☆ [Cited by 459](#) [Related articles](#) [All 9 versions](#)

# Criticism

- Experiments section is not detailed:
  - Software and hardware description is shallow.
  - Why did they use the held-out data in the training?
  - Why did they not use the same amount of features for each variant during the experiment?

# Outline

1. Introduction
2. Log-Linear Models in NLP and L1
3. SGD - L1
4. OWL-QN
5. Experiments
6. Criticism
7. **Conclusions**

# Conclusions

- As future work, they discussed that it would be interesting to tune the metaparameters of the learning rate in an automatic way. Existing algorithms are too expensive.
- They presented a variant of SGD that can efficiently train L1-regularized log-linear models.
- The authors demonstrated empirically that their training algorithm can produce compact and accurate models much more quickly than OWL-QN.
- Trying to find the shortest/straight path to a minimum is expensive and complex.

# Annex 1

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$

$$= f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$$

image source: Stewart,  
2011

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \cdot J(x_1, x_2, x_3, \dots, x_n) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} & \dots & \frac{\partial f_2}{\partial x_n} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} & \dots & \frac{\partial f_3}{\partial x_n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \frac{\partial f_n}{\partial x_3} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

image source:  
brilliant.org

image source:  
chegg.com

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$f'(a) = \lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$$

image source: Stewart,  
2011

# Annex 2

A vector  $\mathbf{d}$  is referred to as a descent direction at  $x$  if  $f'(x; \mathbf{d}) < 0$

$$f'(x; \mathbf{d}) = \lim_{\alpha \downarrow 0} \frac{f(x + \alpha \mathbf{d}) - f(x)}{\alpha}$$

Sign function  $\sigma$

$$\sigma(x) = \begin{cases} -1 & x < 0 \\ 0 & x = 0 \\ 1 & x > 0 \end{cases}$$

The projection of  $x$  onto an orthant defined by  $y$

$$\pi(x; y) = \begin{cases} x_i & \text{if } \sigma(x_i) = \sigma(y_i) \\ 0 & \text{otherwise} \end{cases}$$

if  $\mathbf{B}_k$  is the approximated Hessian matrix of a smooth function  $f$  at  $x^k$ , and  $g^k$  is the gradient of  $f$  at  $x^k$ :

$$Q(x) = f(x^k) + (x - x^k)^T g^k + \frac{1}{2}(x - x^k)^T \mathbf{B}_k (x - x^k)$$

the value that minimizes  $Q$  is computed according to  $x^* = x^k - \mathbf{H}_k g^k$  where  $\mathbf{H}_k = \mathbf{B}_k^{-1}$   
a Quasi Newton Method explores

$$x^{k+1} = x^k - \alpha \mathbf{H}_k g^k, \text{ for } \alpha \in (0, \infty)$$

(Andrew,  
2007)

and maintains only the curvature information from the most recent points

$$s^k = x^k - x^{k-1}$$

$$y^k = g^k - g^{k-1}$$

Choosing an orthant,  
pseudo-gradient

$$\diamond_i f(x) = \begin{cases} \partial_i^- f(x) & \text{if } \partial_i^- f(x) > 0 \\ \partial_i^+ f(x) & \text{if } \partial_i^+ f(x) < 0 \\ 0 & \text{otherwise} \end{cases}$$

and

$$\partial_i^\pm f(x) = \frac{\partial}{\partial x_i} \ell(x) + \begin{cases} C\sigma(x_i) & \text{if } x_i \neq 0 \\ \pm C & \text{if } x_i = 0 \end{cases}$$

The orthant to explore is the one that contains  $x^k$  into which  $-\diamond f(x^k)$  leads:

$$\xi_i^k = \begin{cases} \sigma(x_i^k) & \text{if } x_i^k \neq 0 \\ \sigma(-\diamond f(x^k)) & \text{if } x_i^k = 0 \end{cases}$$

Constraint linear search

$$x^{k+1} = \pi(x^k + \alpha p^k; \xi^k)$$

For any sign vector  $\xi \in \{-1, 0, 1\}^n$ , let us define

$$\Omega_\xi = \{x \in \mathbb{R}^n : \pi(x; \xi) = x\},$$

which is the intersection of an orthant and a plane constraining some coordinates to be zero. Clearly, for all  $x$  in  $\Omega_\xi$ ,

# Curriculum Learning

Seminar Data Analytics 1

International Master's Program in Data Analytics  
University of Hildesheim  
Summer Semester 2018

Shabanaz Chamurally

# Outline

1. General Information about the paper
2. What is the paper about?
  - Motivation
  - Objectives
3. Contributions
4. A curriculum as a continuation method
  - Continuation method
  - How a curriculum can act as a continuation method
5. Methodology
6. Experiments
7. Potential advantages of Curriculum Learning
8. Relation to other machine learning approaches
9. Conclusion

# 1. General Information about the paper

- **Authors:** Yoshua Bengio, Jérôme Louradour, Ronan Collobert, Jason Weston
- **The paper:** Published in *Proceedings of the 26<sup>th</sup> International Conference on Machine Learning*, Montreal, Canada, 2009.
- **Citation:** Cited 748 times (as of 4/5/2018)
- **Available at:**  
<https://qmro.qmul.ac.uk/xmlui/bitstream/handle/123456789/15972/Bengio%20%202009%20Curriculum%20Learning.pdf?sequence=1>

## 2. What is the paper about?

### - *Motivation*

“Humans and animals learn much better when the examples are not randomly presented but organized in a meaningful order ... Training is highly organized based on an education system and a curriculum which introduces different concepts at different times ... to ease the learning of new abstractions.”

From previous research (Elmann, 1993; Rohde & Plaut, 1999; Krueger & Dayan, 2009) which have raised the question:

**Can machine learning algorithms benefit from similar training strategy?**

## 2. What is the paper about?

### - *Objectives*

- Formalize curriculum learning strategies in the context of machine learning
- Explore curriculum learning in various setups
- Show that significant improvements in generalization can be achieved
- Hypothesis that curriculum learning has an effect on the speed of convergence
- In the case of non-criteria, curriculum learning can affect the quality of the local minima
- How a curriculum can be seen as a continuation method

# 3. Contributions

- Explore cases involving vision and language tasks that show that a curriculum learning can benefits machine learning
- Show how curriculum strategies improved generalization and give rise to faster convergence
- Introduce a hypothesis which may help to explain why curriculum learning can help find a better local minima of a non-convex training criterion
- Explore the relation of curriculum learning with other strategies proposed

## 4. A curriculum as a continuation method - Continuation Methods

- Continuation methods are optimization strategies for dealing with minimizing non-convex criteria.
- Involves a sequence of training
  1. Starting from one that is easier to optimize
  2. Increasing difficulty in optimization
  3. Ending with training criterion of interest

## 4. A curriculum as a continuation method - Continuation Methods

- E.g.

<b>Cost Function: <math>C_\lambda(\theta)</math></b>	
<b>Objective</b>	<b>Minimize <math>C_1</math></b>
<b>Criteria easily optimized</b>	<b><math>C_0</math></b>

- Using continuation methods:
  1. Minimize  $C_0$
  2. Gradually increase  $\lambda$  while keeping  $\theta$  at a local minimum
  3. Till reach the criterion of interest  $C_1$

## 4. A curriculum as a continuation method

### - How a curriculum can be a continuous method

- Introduce the idea that a curriculum can also be a sequence of training criteria.
- Each training criteria is associated with a **different set of weights** on the training examples.
  1. Initially, the **weights** favour “**easier**” examples or examples illustrating the **simplest** concepts that can be learned most easily.
  2. The next training criterion is a **slight change** in the **weighting** of example.
  3. The **probability** of sampling **slightly more difficult** examples **increases**.
  4. At the **end** of the sequence, the **reweighting** of the examples is **uniform**.

# 5. Methodology

## - *Notations*

<b>Notations</b>	<b>Description</b>
$\mathcal{Z}$	Random variable representing an example for the learner
$P(\mathcal{Z})$	Target training distribution from which the learner should learn a function
$0 \leq \lambda \leq 1$	Training step
$0 \leq W_\lambda \leq 1$	Weight applied to example $\mathcal{z}$ at step $\lambda$
$W_1(\mathcal{z}) = 1$	Weight applied to $\mathcal{z}$ at the end ( $\lambda = 1$ )
$Q_\lambda(\mathcal{z})$	Training Distribution at step $\lambda$
$H(Q_\lambda)$	Entropy of training distribution at step $\lambda$

# 5. Methodology

## - *Description*

- The training distribution at step  $\lambda$  is:

$$Q_\lambda(\mathbf{z}) \propto W_\lambda(\mathbf{z}) P(\mathbf{z}) \quad \forall \mathbf{z}$$

- Such that  $\int Q_\lambda(\mathbf{z}) d\mathbf{z} = 1 \quad \forall \mathbf{z}$
- By monotonically increasing  $\lambda$ , from  $\lambda = 0$  to  $\lambda = 1$ 
  - $Q_\lambda$  is called a curriculum if the *entropy of these distributions increases*

$$H(Q_\lambda) < H(Q_{\lambda + \epsilon}) \quad \forall \epsilon > 0$$

- $W_\lambda(\mathbf{z})$  is also monotonically increasing

$$W_{\lambda + \epsilon}(\mathbf{z}) \geq W_\lambda(\mathbf{z}) \quad \forall \mathbf{z}, \forall \epsilon > 0$$

# 5. Methodology

## - *Illustration*

1.  $Q_\lambda$  is a **finite set** of examples
2. Increasing  $\lambda$  means adding **new** examples to that set
3. Support of  $Q_\lambda$  **increases** with  $\lambda$ 
  - The **entropy increase** so as to **increase** the **diversity** of training examples
4. The training distributions:
  1. **Start** with a small set of **easy** examples
  2. **Slight** change in  $W_\lambda(\mathbf{z})$  that **increase** the probability of more **difficult** examples as they get added.
  3. At the **end** all examples will have an **equal** weight of 1.
  4. **Train** on the target training set.

## 6. Experiments: *Toy Experiments with a convex Criterion* - *Introducing Gradually More Difficult Examples speeds-up Online Training*

Train a Perceptron from artificially generated data by introducing gradually more difficult examples.

### Notations:

<b>Target: <math>y = \text{sign}(w'x_{\text{relevant}})</math></b>	
<b><math>w</math></b>	Weight vector of a Bayes classifier
<b><math>(x, y)</math></b>	Training pairs
<b><math>x_{\text{relevant}}</math></b>	Relevant input sampled from a Uniform(0,1)
<b><math>x_{\text{irrelevant}}</math></b>	Irrelevant input, not predictive of the target class. Either set to 0 or sampled from a Uniform (0,1)

## 6. Experiments: *Toy Experiments with a convex Criterion* - *Introducing Gradually More Difficult Examples speeds-up Online Training*

### Definition of Easy and Difficult examples

- Train the perceptron from **easy** examples to most **difficult** ones.
  - Easy examples: with **all** irrelevant inputs **zeroed out**
  - Difficult examples: with **none** of the irrelevant inputs **zeroed out**.
- Another way to differentiate between easy and difficult
  - By the margin  $yw'x$ .
  - Easy examples: **larger** margin
  - Difficult examples: **smaller** margin values
- Increase the difficulty levels in each step by increasing the weights of the examples.

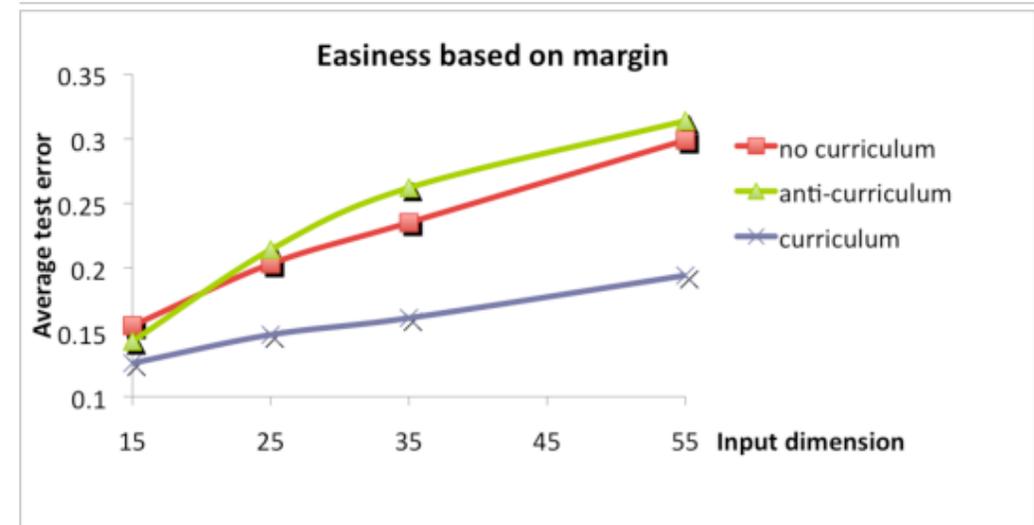
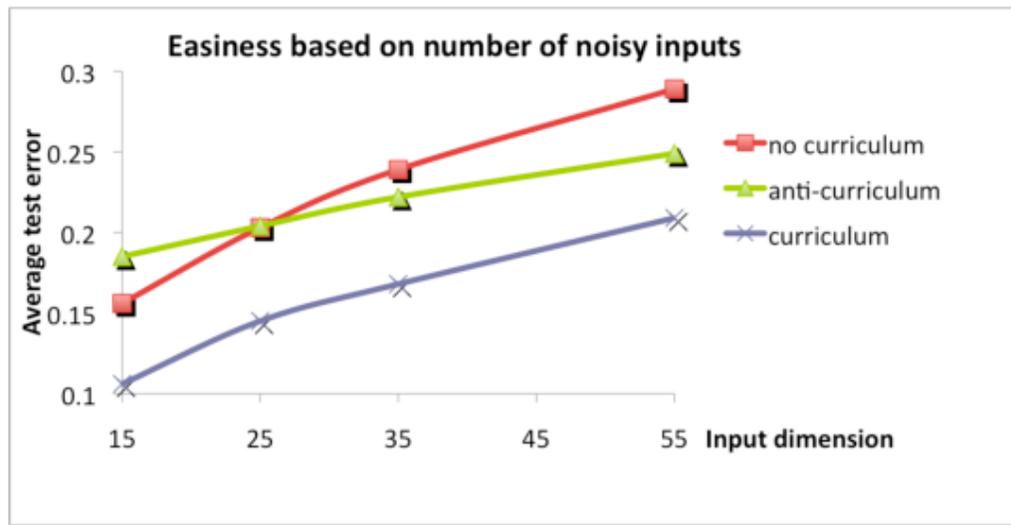
## 6. Experiments: *Toy Experiments with a convex Criterion* *- Introducing Gradually More Difficult Examples speeds-up Online Training*

### **The Experiment**

- Train the perceptron with 200 examples
  - i.e. 200 perceptron updates
- Measure generalization error at the end of training
- Averaged error across 500 repetitions
  - from different initial conditions
  - and different random sampling of training examples

## 6. Experiments: *Toy Experiments with a convex Criterion* - *Introducing Gradually More Difficult Examples speeds-up Online Training* *Online Training*

### Results



- Error rate between the curriculum strategy and the no-curriculum strategy are statistically significant.

# 6. Experiments: On Shape Recognition

## Objective

- Classify geometrical shapes into 3 classes:
  - Rectangle
  - Ellipse
  - Triangle

## Input

- 32 x 32 grey-scale image

## Datasets

- BasicShapes
- GeomShapes

# 6. Experiments: On Shape Recognition

## Definition of Easy examples

### ➤ BasicShapes:

- Squares
- Circles
- Equilateral triangles



# 6. Experiments: On Shape Recognition

## Definition of Difficult examples

### ➤ GeomShapes:

- Rectangles
- Ellipses
- Triangles



## 6. Experiments: On Shape Recognition

### **Difference between BasicShapes and Geomshapes:**

- BasicShapes exhibit less variability in shape

### **Similarities between BasicShapes and Geomshapes:**

- Variabilities in both:
  1. Object position
  2. Size
  3. Orientation
  4. Grey levels of the foreground and background

# 6. Experiments: On Shape Recognition

## The Experiment

- Carried out on a multi-layer neural network with 3 hidden layers
- The curriculum consists in a 2-step schedule:
  1. Perform gradient descent on the BasiShapes training set, until “switch epoch” is reached
  2. Then perform gradient descent on the GeomShapes training set.
- An epoch is a stochastic gradient descent pass through a training set of 10000 examples.

# 6. Experiments: On Shape Recognition

## The Experiment

- After switching to the target training distribution, training continues until:
  1. 256 epochs
  2. Or validation set error reaches a minimum (early stopping)
- Used 10000 examples in both sets, 5000 examples for validation and 5000 for testing.

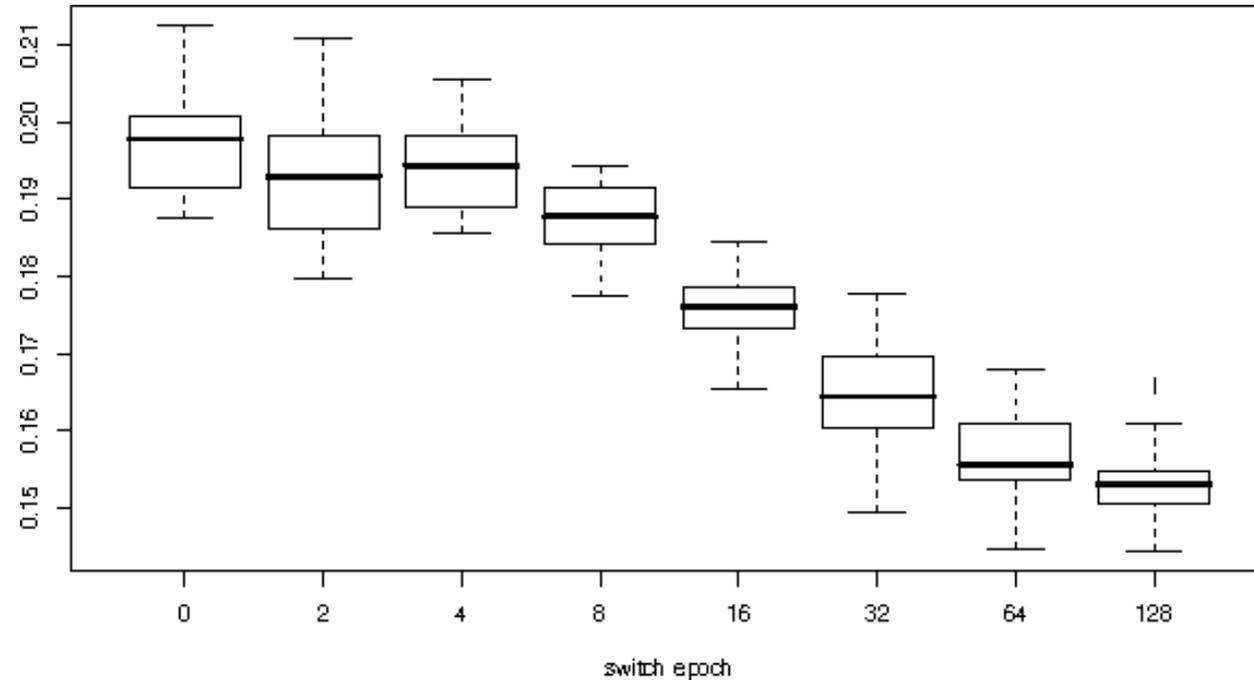
# 6. Experiments: On Shape Recognition

## Results

- Generalization error is evaluated on GeomShapes test set
- **Baseline:** Training only on GeomShapes training set (for the same number of training epochs).  
**Switch epochs = 0**
- Distribution test errors:
  1. Over 20 different random seeds
  2. Different values of switch epochs
    - 0: The baseline with no curriculum
    - Powers of 2 until 128

# 6. Experiments: On Shape Recognition

## Results



- Figure shows the distribution test error after early stopping.
- Each box corresponds to 20 seeds for initializing the parameters.
- Horizontal line: Median (50<sup>th</sup> percentile)
- Borders: 25<sup>th</sup> and 75<sup>th</sup> percentile
- End of bars: 5<sup>th</sup> and 95<sup>th</sup> percentile

### Observation:

- Best generalization is obtained by doing a 2 stage curriculum.

# 6. Experiments: On Shape Recognition

## Potential issue

Curriculum-trained model overall saw more examples than the no-curriculum examples.

To eliminate the explanation that better results are obtained with the curriculum-trained model:

- Trained a no-curriculum with union of BasicShapes and GeomShapes:
  - Final test error significantly worse than with the curriculum
- Trained only with BasicShapes:
  - Yielded poor results

# 6. Experiments: On Language Modeling

## Objective

- Train a language model for predicting the best word which can follow a given context of words in a correct English sentence.
- Compute a score for the next word that will have a large rank compared to the scores of other words
  - Given any fixed size window of text  $s$ , a language model  $f(s)$  produces a score for these windows of text.
  - The score of a correct windows  $s$  should be larger, by a margin of 1, than any other word sequence  $s^w$  where the last word has been replaced by another word  $w$ .

# 6. Experiments: On Language Modeling

## Ranking Loss Function

The ranking loss over sequences  $\mathbf{s}$  sampled from a dataset  $\mathcal{S}$  of valid English text windows:

$$C_{\mathcal{S}} = \sum_{w \in \mathcal{D}} \frac{1}{|\mathcal{D}|} C_{\mathcal{S},w} = \sum_{w \in \mathcal{D}} \frac{1}{|\mathcal{D}|} \max(0, 1 - f(\mathbf{s}) + f(\mathbf{s}^w))$$

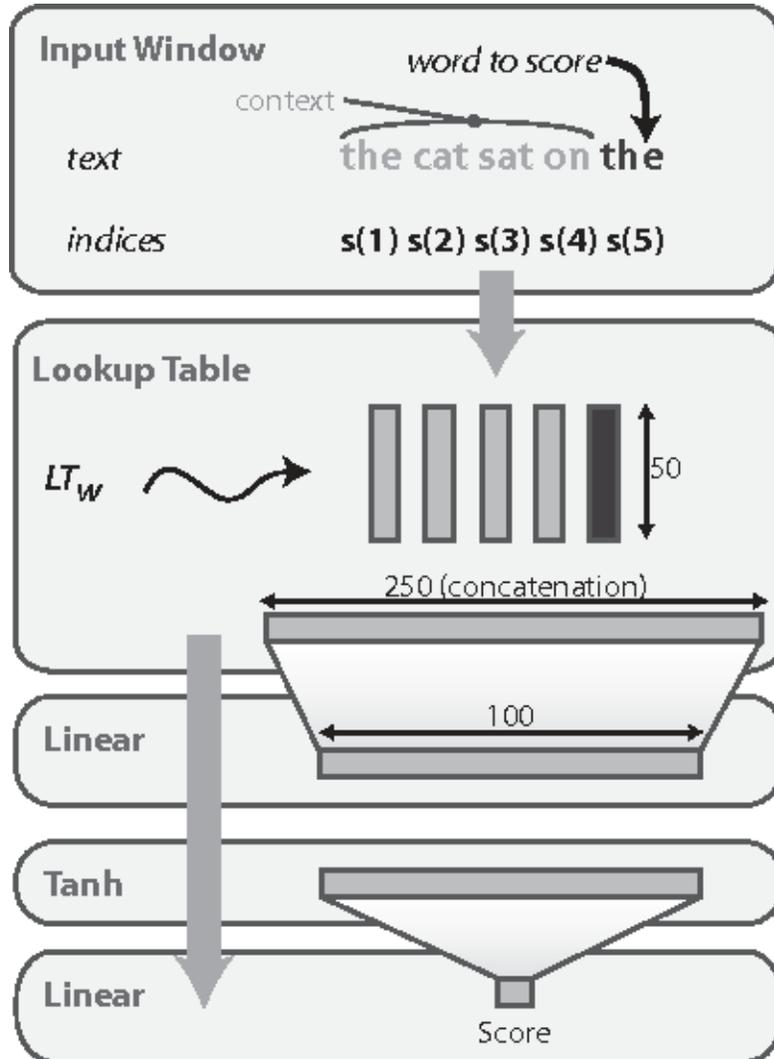
Where  $\mathcal{D}$  is the word vocabulary.

For each word sequence  $\mathbf{s}$

- $f(\mathbf{s})$  and  $f(\mathbf{s}^w)$  is computed
- And the gradient of  $\max(0, 1 - f(\mathbf{s}) + f(\mathbf{s}^w))$

# 6. Experiments: On Language Modeling

## Architecture of language model



- The first layer extracts features for each word in a sentence. Map words into real-valued vectors
- The second layer extracts features from the sentence treating it as a sequence with local and global structure (i.e., it is not treated like a bag of words).
  - Each word  $i$  is embedded into a  $d$ -dimensional space using a look-up table  $LT_w(.)$
  - An input window  $\{s_1, s_2, \dots, s_n\}$  of  $n$  words in  $D$  is transformed into a series of vectors  $\{W_{s_1}, W_{s_2}, \dots, W_{s_n}\}$  by applying the look-up table to each of its words
- Vectors Obtained by the look-up table layer are then concatenated and fed to a classical linear layer.
- Following the 3 above layers is a non-linearity,  $\tanh(.)$
- The score of the language model is finally obtained after a linear layer.

# 6. Experiments: On Language Modeling

## The Experiment

- Training set  $S$  is all the possible windows of text of size = 5 from Wikipedia
- Obtained 631 million of windows
- The curriculum is chosen to grow the vocabulary size:
  1. **First pass** over Wikipedia was performed using the **5000 most frequent** words in the vocabulary
  2. Which is then **increased by 5000 words** at each subsequent pass through Wikipedia
  3. At each pass, any window of text containing a word **not** in the considered dictionary was **discarded**.
- The training set is thus increased after each pass through Wikipedia.

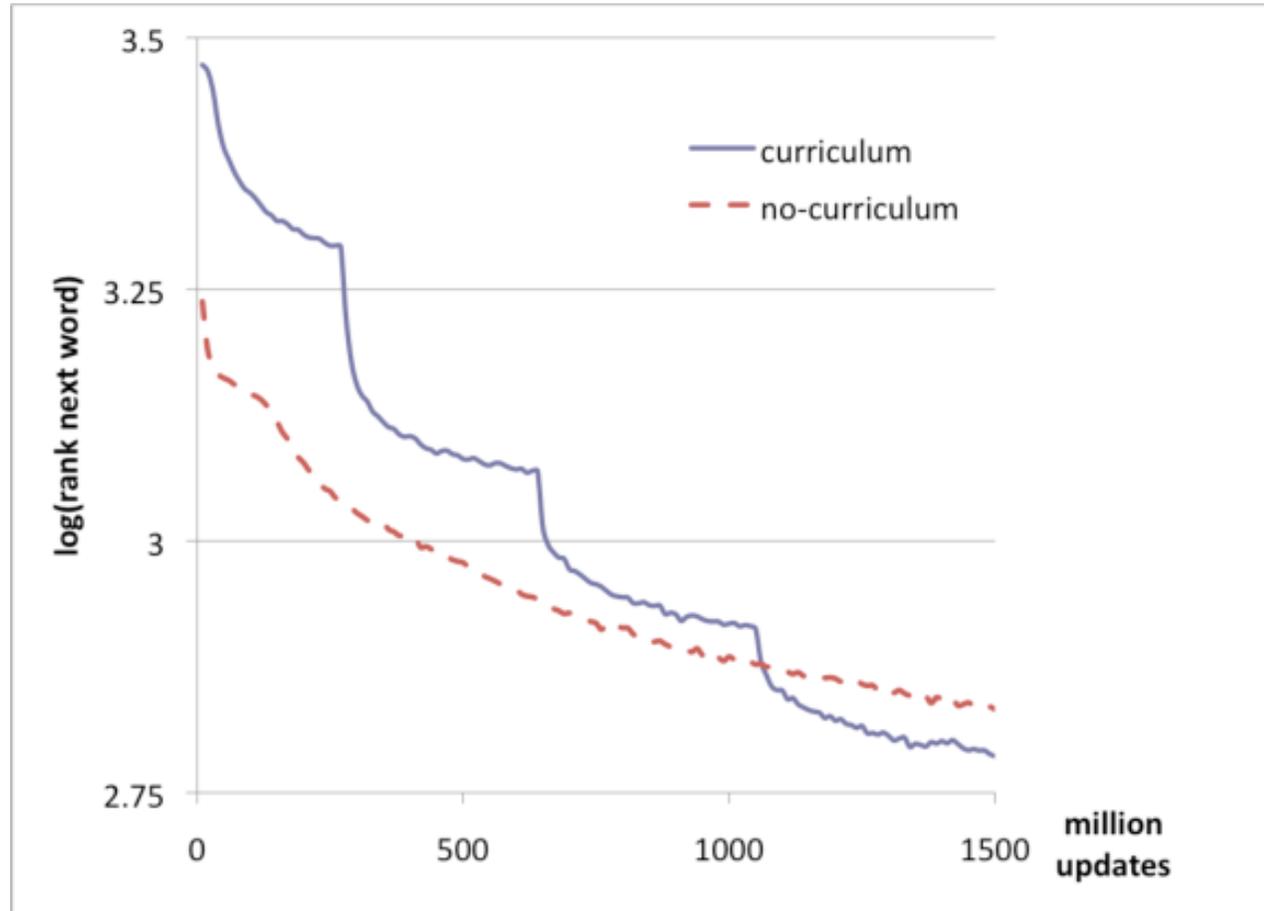
# 6. Experiments: On Language Modeling

## Results

- The ranking loss is **minimized** using stochastic gradient descent.
- Compare against no curriculum
  - Where the network is trained using **final desired vocabulary** size of 20,000.
- Evaluation criterion:
  - Is the average of the **log of the rank** of the last word in each test window
  - Taken in a test set of 10,000 windows not seen during training
  - With words from the most 20,000 frequent ones.

# 6. Experiments: On Language Modeling

## Results



- The log rank on the target distribution with the curriculum strategy **crosses the error** of the no-curriculum strategy after about 1 billion updates.
  - Shortly after switching to the target vocabulary size of 20 000 words and the difference keeps increasing afterwards.
- The final test set average log-ranks are 2.78 and 2.83 respectively and the difference is **statistically significant**.

## 7. Potential advantages of a curriculum strategy

- **Faster training time** in the online setting because the learner **wastes less time** with noisy or harder to predict examples, when the learner is not ready yet.
- Guiding training towards **better local minima** in parameter space, specifically useful for non-convex methods.
- A curriculum can be seen as a **continuation method**

## 8. Relation to other machine learning approaches

- **Unsupervised pre-training procedures:** An appropriate curriculum strategy acts to **help** the training process
  - **faster convergence** to better solutions and both methods have a **regularizing effect** and **lower** the generalization error for the same training error
- **Active Learning:** The learner **benefits most** from focussing on “interesting examples” which are close to the learner's frontier of knowledge instead of focussing on examples that are “too easy” or “too hard”. The learner could succeed to capture and gradually expand its border

## 8. Relation to other machine learning approaches

- **Boosting algorithm:** **Difficult** algorithm are **gradually** emphasized.
  - However, a curriculum starts by focussing on the easier examples rather than a uniform distribution over trained set
- **Transfer and lifelong learning:** Initial tasks are used to **guide** the learner so that it will **perform better** on the final task learning to improve generalization.
  - Curriculum learning **adds** the notion of **guiding** the optimisation process to **converge better** or towards better local minima

# 9. Conclusion

- Started with the question “**Can machine learning algorithm benefit from a curriculum strategy?**”
- From the experiments, it is **plausible** that some curriculum strategies work better than others
- Better results could be obtained with more appropriate curriculum strategies
- It would be useful to understand general principles that makes some curriculum work better than others and why.

# Criticism

- In the experiments concerning the convex criterion:
  - Anti-curriculum performs better than no-curriculum.
  - Given that curriculum learning is modelled on the idea that learning benefits when examples are presented in order of increasing difficulty, one would expect anti-curriculum to perform worse.
- The results of the experiments were all compared to the results of a no-curriculum model:
  - The paper could have explored how a curriculum strategy can perform better than existing machine learning approaches.

# Combined Regression & Ranking

Presentation on the paper by D. Sculley (2010)

Daniela Thyssens

# Outline

1. General information about the paper
  2. What is the paper about?
  3. Methodology
    - Learning in general & Notation
    - Minimizing Loss & Loss Functions
    - Overall Optimization Problem & Algorithm
    - Learning with Stochastic Gradient Descent
  4. Experiments
    - Setup and Results for each experiment
  5. Future work
  6. Conclusion and own thoughts
- Appendix

# 1. General information about the paper

- **The author:** D. Sculley is a Manager, Tech Lead and Staff Software Engineer at Google, Inc. (Research at Google)
- **The paper:** Published in the Proceedings of the 16<sup>th</sup> ACM SIGKDD international conference
  - The conference is hosted by SIGKDD (specialised interest group on knowledge discovery and data mining)

# What is the paper about?

- The paper argues that combining regression and ranking methods is beneficial for many real world tasks.
  - General throwback: Doesn't good regressing performance already include accurate ranking?
  - On the other hand one could ask why not (combining both)?
- What is fundamental about it? (for text analysis and application)
  - Typical supervised machine learning problem
  - Easy to implement algorithm – stochastic gradient descent (popular method)

# What is the paper about?

## - *Goal of the paper*

- Claim:

The combined method achieves the same performance as (& sometimes outperforms) regression- and ranking-only methods for large scale experiments

- The combined regression and ranking (CRR) algorithm gives optimized regression and ranking results simultaneously
- CRR as an “easy to implement” – optimization algorithm
  - Through the use of stochastic gradient descent
- Efficient for the use of large-scale data sets

# What is the paper about?

## - *Motivation: Ranking vs. Regression*

- Case: extreme minority classes
  - example, the query (=target) for a search on the web is very specific.
  - Disproportionately high amount of target values  $y=0$  → easy to predict
  - But almost useless trying to optimize ranking
    - Because features of minority class ( $y=1$ ) are treated as noise and hence mainly ignored, there's a high probability of them being ranked "not relevant"
- we still need good regression **and** good ranking for web-search results
- Case: order-preserving transformation in ranking
  - Invariance of pairwise loss

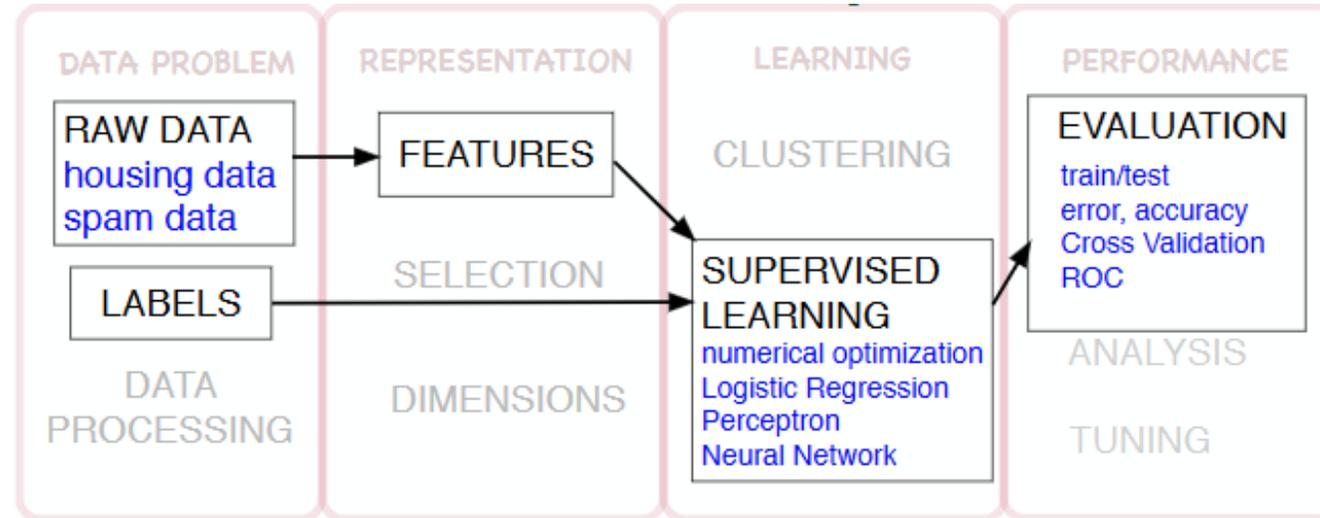
# What is the paper about?

## - *Motivation: Why combining?*

- “Many real world tasks require the achievement (...)” of both
  - (D. Sculley, Combined Regression and Ranking, 2010, p.1)
- Arguments:
  - Perfect regression results do not imply good ranking (extreme minority class)
  - Don't loose out on performance for both goals
  - Overall enhancement of the regression result (rare events)
- Examples of where both ranking and regression are useful:
  - Text classification (Experiment 1)
  - Ranking concerning Web-search results (Experiment 2)
  - Click-Prediction for sponsored ads – commercial application (Experiment 3)

# Methodology

## *-Learning and testing*



- Here: Labels  $\rightarrow y$   
Features  $\rightarrow x$   
Superv. Learning  $\rightarrow$  optimization – minimizing loss functions  
(through stochastic gradient descent)

# Methodology

## -Notation & Preliminaries

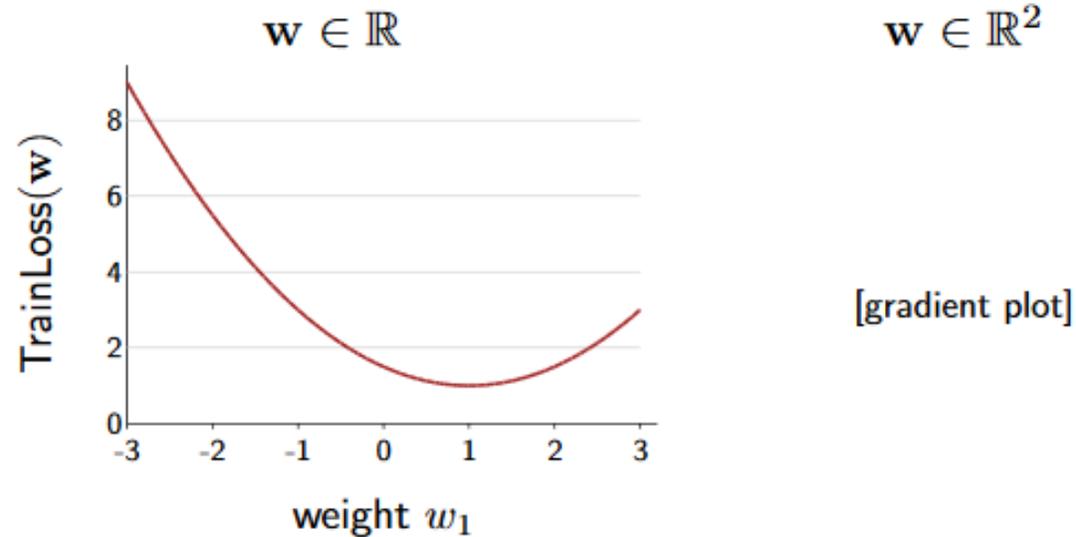
- “Input” space:  $X$  and “Output” space:  $Y$
- Training set  $D$  drawn i.i.d. from  $\mu(z)$  on  $Z = X \times Y$

	Dataset	Variables	Training Examples	Model $w$
Regression	$D$	$x$ : Feature Vector $y$ : Label $q$ : Identifier	Tuples $(x,y,q)$	predicts target value $y$ given $x$ → <i>supervised regression</i> through $f(w,x)$
	expanded version ↓			
Ranking	$D$ (expanded by P candidate pairs)	$a,b$ : Candidate Pairs $y$ : Label $q$ : Identifier	Tuple Pairs $(a,y_a,q_a)$ & $(b,y_b,q_b)$	optimizes the pairwise objective function → <i>supervised ranking</i> $f(w, a - b)$

- $y_a > y_b$  implies  $a$  is preferred to/ranked better than  $b$
- $q$  is the identifier for the query shard – if there is no query shard :  $q = 1$

# Methodology

## *-minimizing Train-Loss*



- Find  $\mathbf{w}$  such that  $\text{TrainLoss}(\mathbf{w})$  is minimized
- Note: the plot is one dimensional, but  $\mathbf{w}$  is actually  $R$ -dimensional  
→ in general, not possible to set  $\mathbf{w}$  to a single value that makes every example have low loss.

# Methodology

## *-Loss functions measuring TrainLoss*

- Supervised regression:

$$L(w, D) = \frac{1}{|D|} \sum_{(x,y,q) \in D} l(y, f(w, x)) \longrightarrow \text{(aggregated loss)}$$

- Supervised ranking:

$$L(w, P) = \frac{1}{|P|} \sum_{((a,y_a,q_a),(b,y_b,q_b)) \in P} l(t(y_a - y_b), f(w, a - b))$$

- Function type: standard convex

- Squared loss  $\rightarrow$  classical LMS **regression** method

$$l(y, y') = (y - y')^2$$

- Logistic loss  $\rightarrow$  predicting real valued probability scores (**ranking**)

$$l(y, y') = y \log(y') + (1 - y) \log(1 - y')$$

# Methodology

## *-Loss function Example*

- Squared loss Example :

- $l(y, f(w, x)) = (y - f(w, x))^2$       Where       $f(w, x) = \mathbf{w} \cdot \mathbf{x}$

- Example:  $\mathbf{w}=[2,-1]$ ,  $\mathbf{x}=[2,0]$  and  $y= -1$

- Then the loss is 25

- $[(-1) - ((2 \times 2) + ((-1) \times 0))]^2 = 25$

# Methodology

## *-Optimization Problem*

- Linearly combining aggregate losses into one min. problem:

$$\min_{\mathbf{w} \in \mathbb{R}} \alpha L(\mathbf{w}, D) + (1 - \alpha)L(\mathbf{w}, P) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\alpha \in [0,1])$$

- (combined) **Empirical Risk Minimization (ERM)** with **L2-regularization**:
  - Linear model represented by  $\mathbf{w}$ : convex data fitting term + “regularizer”

# Methodology

## *-Optimization Problem*

- Linearly combining aggregate losses into one min. problem:

$$\min_{\mathbf{w} \in \mathbb{R}} \alpha L(\mathbf{w}, D) + (1 - \alpha)L(\mathbf{w}, P) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\alpha \in [0,1])$$

- (combined) **Empirical Risk Minimization (ERM)** with L2-regularization
  - Solution should:
    - **Generalize** (convergence with increasing number of examples - predictive)

# Methodology

## *-Optimization Problem*

- Linearly combining aggregate losses into one min. problem:

$$\min_{\mathbf{w} \in \mathbb{R}} \alpha L(\mathbf{w}, D) + (1 - \alpha)L(\mathbf{w}, P) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\alpha \in [0,1])$$

- (combined) **Empirical Risk Minimization (ERM)** with L2-regularization
  - Solution should:
    - **Generalize** (convergence with increasing number of examples - predictive)
    - **Be Stable** (low complexity)

# Methodology

## *-Optimization Problem*

- Linearly combining aggregate losses into one min. problem:

$$\min_{\mathbf{w} \in \mathbb{R}} \alpha L(\mathbf{w}, D) + (1 - \alpha)L(\mathbf{w}, P) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (\alpha \in [0,1])$$

- (combined) **Empirical Risk Minimization (ERM)** with L2-regularization

- Solution should:

- **Generalize (convergence with increasing number of examples - predictive)**

- **Be Stable (low complexity)**

- $\lambda$  is the regularization parameter

- prevents the model from predicting noise

- penalising additional information that doesn't reduce the error of the model significantly

# Methodology

## -CRR Algorithm

- Outputs  $\mathbf{w}$  where the objective achieves min. value
- Sampling from  $P$
- **Stochastic gradient descent** - efficiently reduces training times
  - Stochastic Gradient Step for squared loss:
$$(1 - \eta_i \lambda) \mathbf{w}_{i-1} + \eta_i \mathbf{x} (y - \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle)$$
  - Stochastic Gradient Step for logistic loss :
$$(1 - \eta_i \lambda) \mathbf{w}_{i-1} + \eta_i \mathbf{x} \left( y - \frac{1}{1 + e^{-\langle \mathbf{w}_{i-1}, \mathbf{x} \rangle}} \right)$$

---

### Algorithm 1 Combined Regression and Ranking.

Given: tradeoff parameter  $\alpha$ , regularization parameter  $\lambda$ , training data  $D$ , iterations  $t$

---

```
1:  $\mathbf{w}_0 \leftarrow \mathbf{0}$ 
2: for  $i = 1$  to  $t$  do
3:   pick  $z$  uniformly at random from  $[0, 1]$ 
4:   if  $z < \alpha$  then
5:      $(\mathbf{x}, y, q) \leftarrow \text{RandomExample}(D)$ 
6:   else
7:      $((\mathbf{a}, y_a, q), (\mathbf{b}, y_b, q)) \leftarrow \text{RandomCandidatePair}(P)$ 
8:      $\mathbf{x} \leftarrow (\mathbf{a} - \mathbf{b})$ 
9:      $y \leftarrow t(y_a - y_b)$ 
10:  end if
11:   $\eta_i \leftarrow \frac{1}{i\lambda}$ 
12:   $\mathbf{w}_i \leftarrow \text{StochasticGradientStep}(\mathbf{w}_{i-1}, \mathbf{x}, y, \lambda, \eta_i)$ 
13: end for
14: return  $\mathbf{w}_t$ 
```

---

# Methodology

## *-Learning with Stochastic Gradient Descent*

### Stochastic Gradient Descent

- initialize  $w = 0$
- repeat
  - randomly choose training point  $(x_i, y_i) \in \mathcal{D}_n$
  - $w \leftarrow w - \eta \underbrace{\nabla_w \ell(f_w(x_i), y_i)}_{\text{Grad(Loss on i'th example)}}$
- until stopping criteria met

# Methodology

## *Stochastic Gradient Step*

- Deriving the gradient step for squared loss:  $(1 - \eta_i \lambda) \mathbf{w}_{i-1} + \eta_i \mathbf{x}(y - \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle)$
- The prediction function:  $f(\mathbf{x}, \mathbf{w}) = \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle$ 
  - $\mathbf{w} \leftarrow \mathbf{w}_{i-1} - \eta_i \{-\lambda \mathbf{w}_{i-1} + \nabla_{\mathbf{w}} [(y - f(\mathbf{x}, \mathbf{w}))^2]\}$
  - $\mathbf{w} \leftarrow \mathbf{w}_{i-1} - \eta_i \left\{-\lambda \mathbf{w}_{i-1} + \frac{\partial}{\partial \mathbf{w}} [(y - f(\mathbf{x}, \mathbf{w}))^2]\right\}$
  - $\mathbf{w} \leftarrow \mathbf{w}_{i-1} - \eta_i \{-\lambda \mathbf{w}_{i-1} + \mathbf{x}[y - \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle]\}$
  - $\mathbf{w} \leftarrow (1 - \lambda \eta) \mathbf{w}_{i-1} - \eta_i \{\mathbf{x}[y - \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle]\}$
  - $\mathbf{w} \leftarrow (1 - \eta_i \lambda) \mathbf{w}_{i-1} + \eta_i \mathbf{x}(y - \langle \mathbf{w}_{i-1}, \mathbf{x} \rangle)$
- $i$  denotes the iteration step hence, for each iteration SGD picks **random** training data  $(\mathbf{x}, y)$  and updates  $\mathbf{w}$  on this example only.

# Methodology

## *-Evaluation / Performance Metrics*

- Regression metric:

$$MSE = \frac{1}{|D|} \sum_{(x,y,q) \in D_t} (y - f(\mathbf{w}, \mathbf{x}))^2$$

- Ranking metrics:

$$1 - AUC \text{ Loss}$$

Probability that the model ranks a randomly chosen observation with  $y = 0$  (by mistake) higher than a randomly chosen observation with  $y = 1$

# Methodology

## -Performance Metrics continued

- Ranking metrics:

- **MAP** = arithmetic mean of **AP**(**q**)

$$AP(\mathbf{q}) = \frac{1}{r_q} \sum_{n=1}^{|\mathcal{Q}|} P@n * r(n) \quad \text{where } P@n = \frac{r_n}{n}$$

- Relevant vs. irrelevant examples across query shard set
- **r**(**n**) returns 1 if example **n** is relevant and 0 otherwise.

- **NDCG** =  $Z_n \sum_{j=1}^n \frac{2^{y_j-1}}{\log(1+j)}$

- quality of ranking for multiple levels of relevance

# Experiments

- Benchmark: Checking effectiveness of CRR against natural competitors: regression- and ranking-only.
  - On one hand, a good strategy to prove the “best of both” argument.
  - On the other hand, fairly easy to show that ranking-only method is outperformed by CRR when it comes to tasks where regression takes a considerable role and vice versa for regression-only tasks.
- 2/3 experiments use publicly available benchmark datasets for text-mining (RCV1) and document ranking (LETOR).
- The third experiment, regarding click prediction uses a proprietary dataset.

# Experiments

## *-RCV1 Experiment*

- RCV1 (or RCV1-v2/LYRL2004) is a text categorization test collection distributed as on-line appendices.
  - Large dataset, sparsity, different class distribution levels
  - Use larger split (781,265 examples) of dataset for training, and the smaller (23,149) for testing.
  - Logistic loss
  - 1,000,000 stochastic gradient descent steps
  - 40 topics chosen out of 103
- Goal:
  - predict excellent **ranking** (relevant vs. irrelevant topic identification)
  - and get good overall topic-based relevance scores (**regression** part)

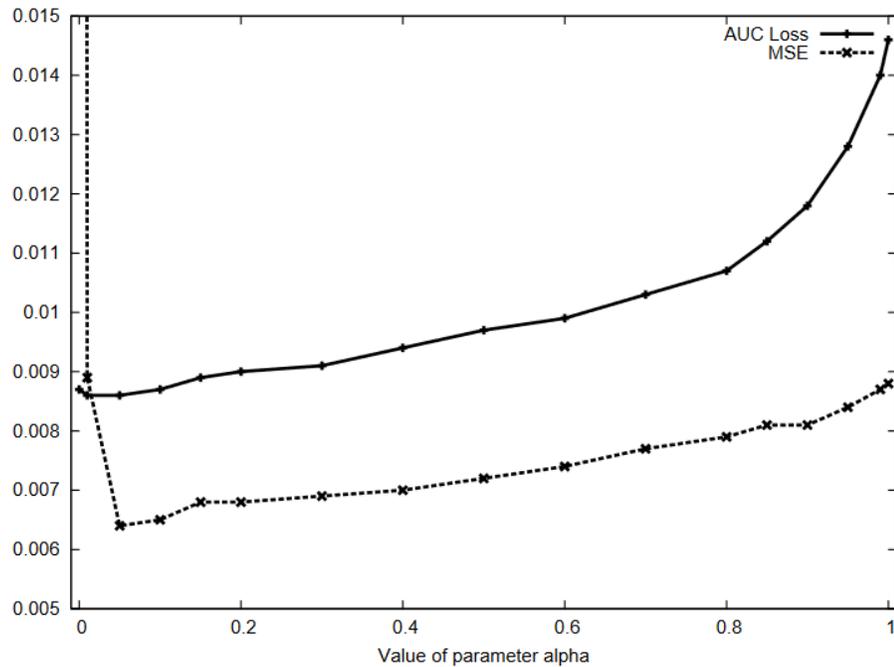
# Experiments

## *-RCV1 Experiment - Results*

- CRR gives desired result: Best of both performance
  - **Best MSE and best AUC Loss** in **16/40** topic-identifications (\*)
  - In 19/40 cases CRR achieves best performance in one metric and in range of 0.001 of the best performance of the other metric.
  - For all **40/40** cases **CRR achieves best performance in at least one metric** (considering that  $\alpha$  is set to 0.5)
  - **CRR performs notably good on the MSE metric concerning minority class distributions (<10%)** and outperforms regression-only in 20/33 cases.
    - adding ranking information to minority class regression problems enhances regression performance.
  - CRR benefits diminish as class-distribution → more evenly balanced

# Experiments

## -RCV1 Experiment -Results



TASK	% POSITIVE	REGRESSION		RANKING		CRR		*
		AUC Loss	MSE	AUC Loss	MSE	AUC Loss	MSE	
E141	0.05%	<b>0.000</b>	0.001	<b>0.000</b>	0.293	<b>0.000</b>	<b>0.000</b>	*
GOBIT	0.06%	0.002	<b>0.001</b>	<b>0.001</b>	0.162	0.002	<b>0.001</b>	*
E61	0.06%	0.002	<b>0.001</b>	<b>0.001</b>	0.320	<b>0.001</b>	<b>0.001</b>	*
GTOUR	0.10%	0.030	<b>0.001</b>	<b>0.005</b>	0.245	<b>0.005</b>	<b>0.001</b>	*
C331	0.13%	0.003	<b>0.001</b>	<b>0.001</b>	0.205	<b>0.001</b>	<b>0.001</b>	*
E143	0.15%	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	0.296	<b>0.001</b>	<b>0.001</b>	*
G152	0.15%	0.005	<b>0.001</b>	<b>0.003</b>	0.239	<b>0.003</b>	<b>0.001</b>	*
G155	0.16%	0.007	0.002	<b>0.004</b>	0.223	<b>0.004</b>	<b>0.001</b>	*
E411	0.17%	<b>0.002</b>	0.002	<b>0.002</b>	0.289	<b>0.002</b>	<b>0.001</b>	*
C313	0.18%	0.047	<b>0.002</b>	<b>0.014</b>	0.281	0.016	<b>0.002</b>	*
E311	0.19%	<b>0.001</b>	0.002	<b>0.001</b>	0.311	<b>0.001</b>	<b>0.001</b>	*
C32	0.19%	0.019	<b>0.002</b>	<b>0.012</b>	0.180	0.013	<b>0.002</b>	*
G157	0.19%	<b>0.001</b>	0.002	<b>0.001</b>	0.254	<b>0.001</b>	<b>0.001</b>	*
C16	0.21%	0.022	<b>0.002</b>	<b>0.012</b>	0.234	0.013	<b>0.002</b>	*
GWELF	0.22%	0.010	<b>0.002</b>	<b>0.005</b>	0.236	0.006	<b>0.002</b>	*
E513	0.23%	0.004	0.002	<b>0.003</b>	0.300	<b>0.003</b>	<b>0.001</b>	*
E14	0.28%	0.008	0.003	<b>0.003</b>	0.281	0.004	<b>0.002</b>	*
C173	0.33%	0.005	0.003	<b>0.004</b>	0.237	<b>0.004</b>	<b>0.002</b>	*
E121	0.41%	0.007	0.004	<b>0.004</b>	0.261	0.005	<b>0.003</b>	*
GENT	0.46%	0.014	<b>0.004</b>	<b>0.008</b>	0.126	<b>0.008</b>	<b>0.004</b>	*
C34	0.52%	0.018	0.005	<b>0.011</b>	0.231	0.012	<b>0.004</b>	*
GHEA	0.85%	0.007	0.008	<b>0.005</b>	0.140	0.006	<b>0.006</b>	*
C183	0.87%	0.013	0.008	<b>0.009</b>	0.275	0.010	<b>0.006</b>	*
GDEF	1.01%	0.015	0.009	<b>0.009</b>	0.208	<b>0.009</b>	<b>0.007</b>	*
C42	1.48%	0.009	0.010	<b>0.006</b>	0.242	0.007	<b>0.008</b>	*
E211	1.76%	0.013	0.011	<b>0.010</b>	0.245	<b>0.010</b>	<b>0.009</b>	*
E51	2.77%	0.025	0.019	<b>0.019</b>	0.280	0.021	<b>0.016</b>	*
M12	3.16%	0.010	0.015	<b>0.008</b>	0.288	0.009	<b>0.014</b>	*
C24	3.98%	0.031	0.027	<b>0.025</b>	0.157	0.026	<b>0.024</b>	*
GDIP	4.34%	0.019	0.023	<b>0.017</b>	0.188	0.018	<b>0.022</b>	*
M13	6.89%	<b>0.007</b>	<b>0.018</b>	<b>0.007</b>	0.221	<b>0.007</b>	<b>0.018</b>	*
GPOL	7.11%	0.021	<b>0.031</b>	<b>0.020</b>	0.175	0.021	<b>0.031</b>	*
C152	8.34%	0.026	0.036	<b>0.023</b>	0.178	0.024	<b>0.035</b>	*
C151	10.22%	0.010	<b>0.024</b>	<b>0.009</b>	0.188	<b>0.009</b>	0.025	*
M14	10.98%	0.005	0.021	<b>0.004</b>	<b>0.115</b>	<b>0.004</b>	0.022	*
ECAT	14.90%	0.033	0.054	<b>0.030</b>	0.188	0.031	<b>0.053</b>	*
C15	18.05%	<b>0.013</b>	<b>0.036</b>	<b>0.013</b>	0.132	<b>0.013</b>	0.037	*
MCAT	25.41%	0.011	<b>0.039</b>	<b>0.010</b>	0.113	<b>0.010</b>	0.043	*
GCAT	30.11%	<b>0.012</b>	<b>0.043</b>	<b>0.012</b>	0.062	<b>0.012</b>	0.046	*
CCAT	46.59%	<b>0.022</b>	<b>0.067</b>	<b>0.022</b>	0.073	<b>0.022</b>	0.070	*

# Experiments

## *-LETOR Experiment*

- Public benchmark data: LETOR 4.0 learning to rank data set
  - 2 tasks in the data set (MQ2007, MQ2008) drawn from TREC challenge data sets in information retrieval
  - Relevance labels (such as star indication): {0,1,2}
  - Querysharded + dense information retrieval features (TF-IDF similarities, PageRank,...)
  - Commonly used for ranking-only purposes.
- Goal: propose relevant documents to the end user
  - Obtain **accurate ranking**
  - **Predicting actual relevance** value of a document (**regression**)
    - Why?
      - Search engines could annotate (query-independent) relevance scores and quality thresholds for documents.

# Experiments

## *-LETOR Experiment - Results*

- CRR does not give the desired result: **Not** Best of Both performance.
  - Performs **as good as ranking-only**
  - **MSE** metric performs **better for CRR compared to ranking-only**, but only approaches the one for regression-only.
- Overall, the performance is not too bad, considering that class distribution is fairly balanced.

# Experiments

## -LETOR Experiment - Results

METHOD	MQ2007			MQ2008		
	MAP	NDCG	MSE	MAP	NDGC	MSE
REGRESSION	0.456	0.492	<b>0.182</b>	0.464	0.476	<b>0.144</b>
RANKSVM (BASELINE [18])	0.464	0.497	-	0.470	0.483	-
RANKING	<b>0.466</b>	<b>0.500</b>	0.498	<b>0.479</b>	<b>0.490</b>	0.850
<b>CRR</b>	0.465	0.498	<b>0.214</b>	<b>0.479</b>	0.489	0.229

# Experiments (commercial application)

## *-Click Prediction Experiment*

- Advertisement Impression - Data from live traffic on Google Search
  - Several million advertisement impressions
  - Each example: high-dim. **Representation of ad** together **with label 1 (clicked) or 0 (not clicked)**
- Goal: which advertisement to show when user enters a search query?
  - **Usefulness of CRR:**
    - Real-time auctions for ad-placements are ranked based on:  $\text{bid} * (\text{price}(\text{CTR}))$  (**ranking**)
    - ad-pricing is governed by “next price auction” (need to predict the price at the next lower rank) (**regression**)

# Experiments

## - *Click Prediction Experiment - Results*

- CRR gives desired result: **Best of both performance**
- AUC Loss is **equal to that in ranking-only** and MSE **performs as well as in regression-only**
- Improvement in AUC Loss: “small”, but statistically significant

METHOD	AdSet1	
	AUC Loss	MSE
REGRESSION	0.133	<b>0.084</b>
RANKING	<b>0.132</b>	0.094
<b>CRR</b>	<b>0.132</b>	<b>0.084</b>

# Future Work

- Could extend the CRR algorithm to non-linear functions
  - Instead of predictions depending on dot product ( $\langle \mathbf{w}, \mathbf{x} \rangle$ ), one could use similarity measure and map features ( $\mathbf{x}$ ) to new features ( $\mathbf{x}'$ ).
- Applying L1 regularization (sparse models)
- Increase sophistication of ranking functions
  - from pairwise approach to direct optimization of MAP/NDCG
- Verifying CRR performance for predicting bounce rates

# Conclusion & own Thoughts (Results)

- From the experiments CRR gives **2/3** times the **desired result**:  
“Best of Both Performance” (text-mining tasks and click-prediction)
- Concerning text mining, combining regression and ranking enhances **mainly** the results for **minority class distributions**
  - In this case, regression results benefit more from the fact that ranking is added as another information constrain to the optimization problem in rare events.
- Ranking tasks itself are not necessarily improved by regression and hence **CRR does not necessarily outperform ranking-only** (2<sup>nd</sup> experiment)
  - On the other hand, the distribution was predominantly evenly distributed.

# Conclusion & own Thoughts (Methodology)

- The choice of the loss function: is squared loss appropriate?

- We could also use the absolute value, or L1 loss:

$$V(f(x), y) = |f(x) - y|$$

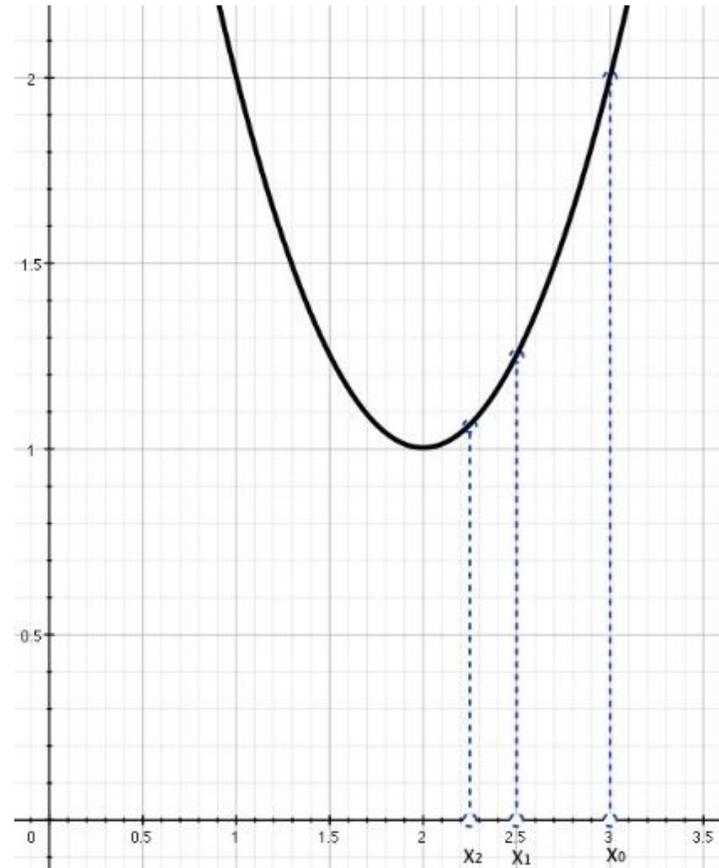
which is less prone to outliers.

- Gradient Descent is usually slow
  - each iteration needs to go over the whole training example set (expensive)

# Appendix-SGD

- Numerical example – Gradient Step
- Suppose:  $J(w) = (w - 2)^2 + 1$ , initial guess for a min:  $w_0 = 3$

- Iteration 1:  $w_1 = w_0 - \lambda \frac{\partial J(w_0)}{\partial w}$   
 $= 3 - 0.25(2 * 3 - 4)$   
 $= 2.5$
- Iteration 2:  $w_2 = w_1 - \lambda \frac{\partial J(w_1)}{\partial w}$   
 $= 2.5 - 0.25(2 * 2.5 - 4)$   
 $= 2.25$
- Iteration 3:  $w_3 = w_2 - \lambda \frac{\partial J(w_2)}{\partial w}$   
 $= 2.25 - 0.25(2 * 2.25 - 4)$   
 $= 2.125$



# Appendix-SGD

- SGD in terms of expectations

- Since  $l(\mathbf{w}) = E_{\mathbf{x}}(l(\mathbf{w}, \mathbf{x})) = \int p(\mathbf{x})l(\mathbf{w}, \mathbf{x})d\mathbf{x}$

- $\nabla_{\mathbf{w}}l(\mathbf{w}) = \nabla_{\mathbf{w}}[E_{\mathbf{x}}(l(\mathbf{w}, \mathbf{x}))] = E_{\mathbf{x}}[\nabla_{\mathbf{w}}l(\mathbf{w}, \mathbf{x})]$

- Then

- $\mathbf{w} \leftarrow \mathbf{w}_{i-1} - \eta_i E_{\mathbf{x}}[\nabla_{\mathbf{w}}l(\mathbf{w}, \mathbf{x})] \rightarrow$  we can't compute this

- Need to approximate with sample data:

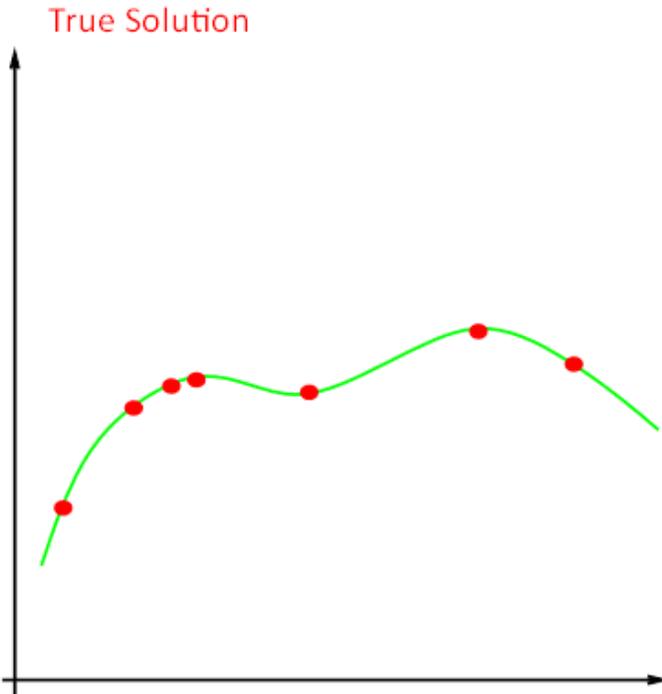
- $\nabla_{\mathbf{w}}l(\mathbf{w}) \approx \sum_{i=1}^D \nabla_{\mathbf{w}}l(\mathbf{w}, \mathbf{x}^i)$

- **So** if we have just one sample (as is the case for SGD- we're updating for each drawn sample), this approximation is very noisy

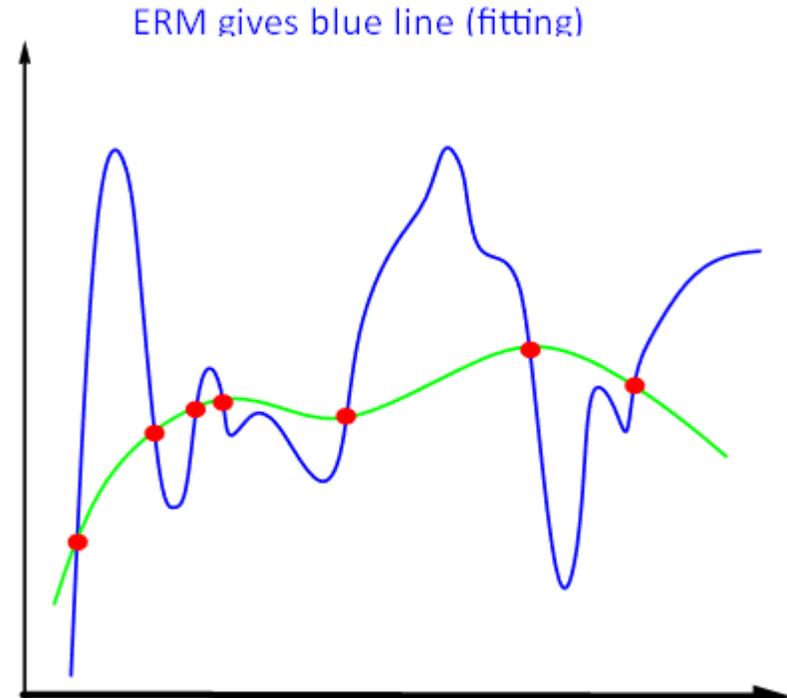
- Need to penalize for this noise  $\rightarrow -\lambda\mathbf{w}_{i-1}$

# Appendix-ERM

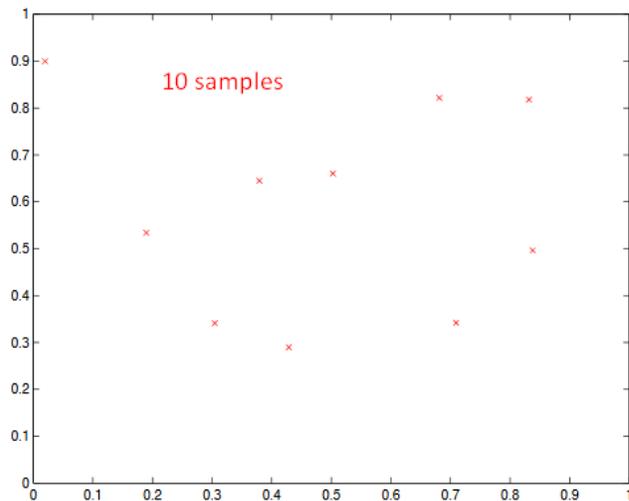
- ERM and generalization



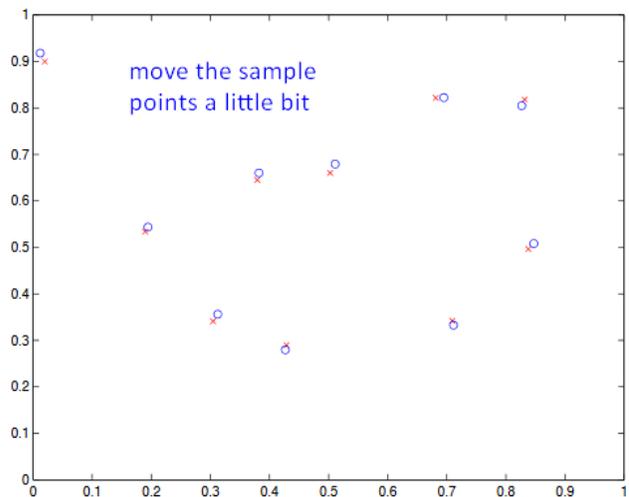
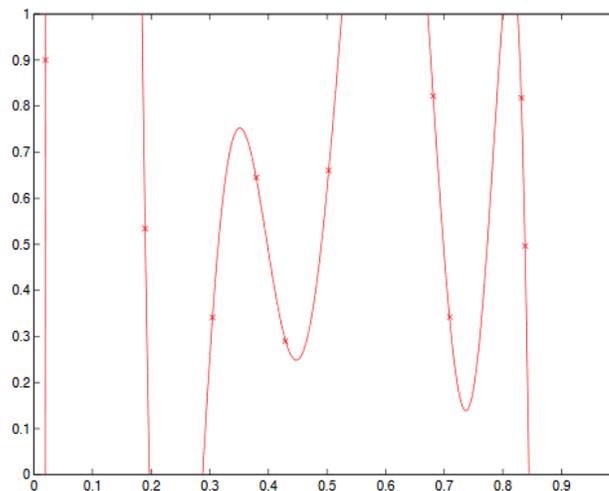
Under which conditions does the ERM line converge with increasing numbers of examples.



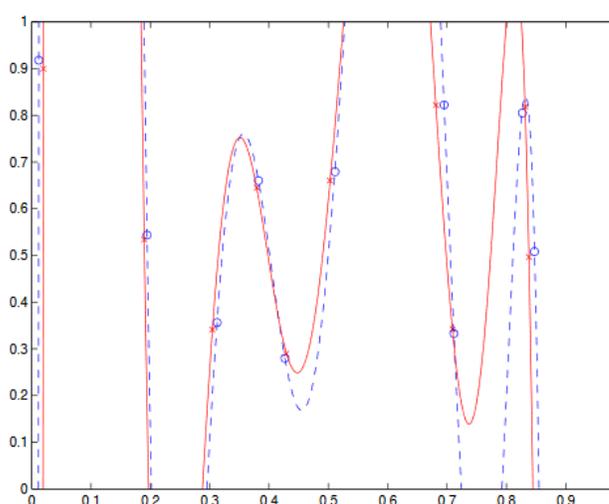
# Appendix-ERM and complexity/stability



→  
fit smoothest  
polynomial  
(high complexity)



→  
model predicts  
relatively bad  
now with this  
fitting  
(the solution  
changes a lot)



# Comparison among the three papers

- Similarities:
  - Training a model
  - Minimizing loss
  - Stochastic gradient descent methods as a fundamental/easy to implement and efficient algorithm step to optimize training time. (practical for large datasets)
- Differences:
  - L1 regularization vs. L2 regularization

# References

- Bengio, Y. (2009). Learning deep architectures for AI. *Foundations & Trends in Mach. Learn.*, to appear
- Bengio, Y., Ducharme, R., & Vincent, P. (2001). A neural probabilistic language model. *Adv. Neural Inf. Proc. Sys.* 13 (pp. 932{938)
- Bottou, L., Bousquet, O. (2008). The Tradeoffs of Large Scale Learning. *In proceedings of the conference "Neural Information Processing Systems 2007". Neural Information Processing Systems 20, 161-168. NIPS.*
- Collobert, R., & Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. *Int. Conf. Mach.Learn. 2008* (pp. 160{167).
- Elman, J. L. (1993). Learning and development in neural networks: The importance of starting small. *Cognition*, 48, 781{799.
- Galen Andrew and Jianfeng Gao. 2007. Scalable training of L1-regularized log-linear models. In *Proceedings of ICML*, pages 33–40.
- Goodfellow, Ian, et al. "Deep Learning." *Deep Learning*, MIT Press, 2016, [www.deeplearningbook.org/](http://www.deeplearningbook.org/).
- Guestrin, C. (2014). CSE546: *Stochastic Gradient Descent*, [Lecture Notes]. Retrieved from <https://pdfs.semanticscholar.org/presentation/767a/4171f9f5c5108caac49400e5871a21b96039.pdf>
- Hsieh, C. (2018). ECS171: *Optimization*, week 4 notes [Lecture Notes]. Retrieved from [http://www.stat.ucdavis.edu/~chohsieh/teaching/ECS171\\_Winter2018/lecture4.pdf](http://www.stat.ucdavis.edu/~chohsieh/teaching/ECS171_Winter2018/lecture4.pdf)
- Krueger, K. A., & Dayan, P. (2009). Flexible shaping: how learning in small steps helps. *Cognition*, 110,380{394.
- Lewis, D. D.; Yang, Y.; Rose, T.; and Li, F. RCV1: A New Benchmark Collection for Text Categorization Research. *Journal of Machine Learning Research*, 5:361-397, 2004. <http://www.jmlr.org/papers/volume5/lewis04a/lewis04a.pdf>.

# References

- Poggio, T. (2011). 9.520: *The Learning Problem and Regularization*, week 2 notes [Lecture Notes]. Retrieved from <http://www.mit.edu/~9.520/spring11/slides/class02.pdf>
- Rosenberg, D. S. (2015). ML1003: *Gradient and Stochastic Gradient Descent* [Lecture Notes]. Retrieved from <https://pdfs.semanticscholar.org/presentation/2c4a/20b112b4160fea7330795853437ba3621733.pdf>
- Rohde, D., & Plaut, D. (1999). Language acquisition in the absence of explicit negative evidence: How important is starting small? *Cognition*, 72, 67-109
- Schwenk, H., & Gauvain, J.-L. (2002). Connectionist language modeling for large vocabulary continuous speech recognition. *International Conference on Acoustics, Speech and Signal Processing* (pp. 765-768). Orlando, Florida
- Sadigh, D. (2018). CS221: *Machine Learning 1*, week 2 notes [Lecture Notes]. Retrieved from <https://web.stanford.edu/class/cs221/lectures/learning1.pdf>
- Sculley, D. (2010, July). Combined regression and ranking. *In Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM.
- Shalev-Shwartz, S., Singer, Y., Srebro, N. (2007). Pegasos: primal estimated sub-gradient solver for SVM. In *ICML '07: Proceedings of the 24<sup>th</sup> international conference on Machine learning*. ICML.
- Stewart, J., 2011. *Calculus: Early Transcendentals*. Cengage Learning.
- Tsuruoka, Y., Tsujii, J. I., & Ananiadou, S. (2009). Stochastic gradient descent training for l1-regularized log-linear models with cumulative penalty. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*
- Yoonkyung, L. (2014). *A statistical View of Ranking: Midway between Classification and Regression* [Lecture Notes]. Retrieved from <http://www.mit.edu/~9.520/spring11/slides/class02.pdf>