# Convergence rate of SGD

- **Theorem**:
  - (see Nemirovski et al '09 from readings)
  - Let $f$ be a strongly convex stochastic function *with param* $\gamma$
  - Assume gradient of $f$ is Lipschitz continuous and bounded

$$\|\nabla f\|_2^2 \leq M^2$$

$$\forall x \quad \|\nabla f(w, x) - \nabla f(w', x)\|_2 \leq L\|w - w'\|_2 \quad L > 0$$

  - Then, for step sizes:

$$\eta_t = \frac{K}{t} \quad K > 0$$

  - The expected loss decreases as O(1/t):

*e.g.* $K = 1/\gamma$

$$\underbrace{E\left[f(w^{(t)}) - f(w^*)\right]}_{\text{how much closer getting to } w^*} \leq \frac{1}{t} L\left(\frac{M^2}{\gamma^2} + \|w^{(0)} - w^*\|_2^2\right)$$

$$\sim O(1/t)$$

24

---

# Convergence rates for gradient descent/ascent versus SGD

$$O(Nd\ln\tfrac{1}{\epsilon}) \longleftrightarrow O\left(\frac{d}{\epsilon}\right)$$

GD --- | --- SGD

- Number of Iterations to get to accuracy

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \epsilon$$

  N data point

  1 data point

  mini batches 100 datapoints

- Gradient descent:
  - If func is strongly convex: O(ln(1/ε)) iterations

- Stochastic gradient descent:
  - If func is strongly convex: O(1/ε) iterations

$$O\left(\ln\tfrac{1}{\epsilon}\right) \text{ iteration} \quad \text{iteration } O(Nd)$$

$$\text{total} = O\left(Nd\ln\tfrac{1}{\epsilon}\right)$$

- Seems exponentially worse, but much more subtle:
  - Total running time, e.g., for logistic regression:
    - Gradient descent:
    - SGD:
    - SGD can win when we have a lot of data

$$O\left(\tfrac{1}{\epsilon}\right) \text{ iterations}, \text{ iteration } O(d)$$

$$\text{total} = O\left(\frac{d}{\epsilon}\right)$$

  - And, when analyzing true error, situation even more subtle… expected running time about the same, see readings

25

1

# Motivating AdaGrad (Duchi, Hazan, Singer 2011)

- Assuming $\mathbf{w} \in \mathbb{R}^d$, standard stochastic (sub)gradient descent updates are of the form:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta g_{t,i}$$

  *(handwritten annotations:)* — step size / — learning rate

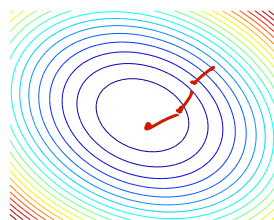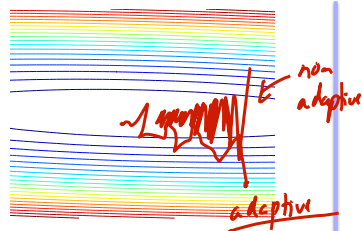  $\eta_{t,i} \leftarrow$ feature-specific step size

  $g_{t,i} \leftarrow$ transforming by a factor $\alpha_i$

- Should all features share the same learning rate?

- Often have high-dimensional feature spaces
  - Many features are irrelevant → *small learning rate*
  - Rare features are often very informative

- Adagrad provides a feature-specific adaptive learning rate by incorporating knowledge of the geometry of past observations

©Emily Fox 2014　　26

---

# Why Adapt to Geometry?

Hard

Nice

*(handwritten: non-adaptive, adaptive)*

| $y_t$ | $x_{t,1}$ | $x_{t,2}$ | $x_{t,3}$ |
|-------|-----------|-----------|-----------|
| 1 | 1 | 0 | 0 |
| -1 | .5 | 0 | 1 |
| 1 | -.5 | 1 | 0 |
| -1 | 0 | 0 | 0 |
| 1 | .5 | 0 | 0 |
| -1 | 1 | 0 | 0 |
| 1 | -1 | 1 | 0 |
| -1 | -.5 | 0 | 1 |

*Examples from Duchi et al. ISMP 2012 slides*

❶ Frequent, irrelevant
❷ Infrequent, predictive
❸ Infrequent, predictive

©Emily Fox 2014　　27

2

# Not All Features are Created Equal
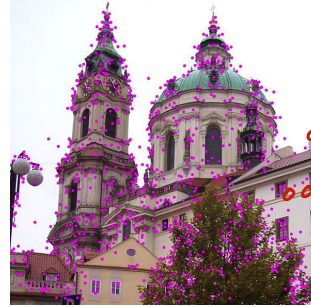
- Examples:

Text data:

The most unsung birthday in American business and technological history this year may be the 50th anniversary of the Xerox 914 photocopier.[a]

---

[a] *The Atlantic*, July/August 2010.

*value word*

High-dimensional image features

*or corners less frequent is more information*

*Images from Duchi et al. ISMP 2012 slides*

---

*Constrained Optimization*          *Original problem*          $\min\limits_{w \in W} \ell(w)$

# Projected Gradient

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta g_{t,i}$$

- Brief aside…   *e.g.,*   $W \Rightarrow \|w\|_1 \leq R$   $\boxed{W}$

- Consider an arbitrary feature space $\mathbf{w} \in \mathcal{W} \subset \mathbb{R}^d$

- If $\mathbf{w} \in \mathcal{W}$, can use ***projected gradient*** for (sub)gradient descent

$$\mathbf{w}^{(t+1)} = \operatorname*{argmin}_{w \in W} \left\| w - \left( w^{(t)} - \eta_t g_t \right) \right\|_2^2 \quad \Leftarrow \quad \text{efficient for some } W$$

*e.g.* $w : \|w\|_2 \leq R$

$\|w\|_1 \leq R$

$\vdots$

$W$          $w^{(t)} - \eta_t g_t$

$w^{(t+1)}$   $w^{(t)}$

*closest point in the space to $w^{(t)} - \eta_t g_t$*

# Regret Minimization

*(handwritten)* $R(T) \to 0$ ⟹ $\frac{1}{T} w^{(t)}, w^{(t+1)} \dots$ as good as $w^*$ — no-regret algorithm

- How do we assess the performance of an online algorithm?

- Algorithm iteratively predicts $\mathbf{w}^{(t)}$ *(handwritten: ad setting, $\hat{y}_t$ click?)*
- Incur **loss** $f_t(\mathbf{w}^{(t)})$ *(handwritten: either click or not)*
- **Regret:**
  What is the total incurred loss of algorithm relative to the best choice of $\mathbf{w}$ that could have been made **retrospectively** *(handwritten: typically $\frac{R(T)}{T} \to 0$ as $T \to \infty$)*

$$R(T) = \sum_{t=1}^{T} f_t(\mathbf{w}^{(t)}) - \inf_{\mathbf{w} \in \mathcal{W}} \sum_{t=1}^{T} f_t(\mathbf{w})$$

*(handwritten: regret; cumulative loss based on sequence of choices; best single $w$ in retrospect; $w^*$)*

---

# Regret Bounds for Standard SGD

- Standard projected gradient stochastic updates: *(handwritten: using $w_{(t)}$)*

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w} \in \mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta g_t)||_2^2$$

- Standard regret bound:

$$\sum_{t=1}^{T} f_t(\mathbf{w}^{(t)}) - f_t(\mathbf{w}^*) \le \frac{1}{2\eta}||\mathbf{w}^{(1)} - \mathbf{w}^*||_2^2 + \frac{\eta}{2}\sum_{t=1}^{T}||g_t||_2^2$$

*(handwritten: $R(T)$; error of $w$ here you started; magnitude of gradients)*

4

# Projected Gradient using Mahalanobis

- Standard projected gradient stochastic updates:

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w}\in\mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta g_t)||_2^2$$

- What if instead of an $L_2$ metric for projection, we considered the **Mahalanobis** norm

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w}\in\mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta A^{-1} g_t)||_A^2$$

*(handwritten annotations)*

$g_t \Rightarrow \begin{pmatrix} g_{t,1} \\ g_{t,2} \end{pmatrix}$  scale up  scale down

$A = \begin{pmatrix} 1 & 0 \\ 0 & 10 \end{pmatrix}$  care more about this gradient

composed by proj. with A

$L_2$ ball: $||w||_2 \le R$    $\sqrt{w^\top w} \le R$

$||w||_A \le R$   $\sqrt{w^\top A w} \le R$   $A = \begin{pmatrix} 1 & 0 \\ 0 & 10 \end{pmatrix}$

$||w||_A \le R$   $A = \begin{pmatrix} \cdot & \cdot \\ \cdot & \cdot \end{pmatrix}$

$A \succeq 0$  positive semi-definite

©Emily Fox 2014    32

---

# Mahalanobis Regret Bounds

*(handwritten top)* in 1d: $||g_t||_{A^{-1}}^2 = \frac{g_t^2}{a}$  min by $a \to \infty$

$tr(A) = \sum_i A_{ii}$

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w}\in\mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta A^{-1} g_t)||_A^2$$

- **What A to choose?**

- Regret bound now:

$$\sum_{t=1}^{T} f_t(\mathbf{w}^{(t)}) - f_t(\mathbf{w}^*) \le \frac{1}{2\eta}||\mathbf{w}^{(1)} - \mathbf{w}^*||_A^2 + \frac{\eta}{2}\sum_{t=1}^{T}||g_t||_{A^{-1}}^2$$

*(handwritten)* $w^{*\top} A w^*$ if $a \to \infty$  $\to \infty$

$||g_t||_{A^{-1}}^2 = g_t^\top A^{-1} g_t = \langle g_t, A^{-1} g_t \rangle$

R(T) want to minimize

- What if we minimize upper bound on regret w.r.t. A in hindsight?

*(handwritten)* Choice of A:  $||g_t||_{A^{-1}}^2$

$$\min_A \sum_{t=1}^{T} \langle g_t, A^{-1} g_t \rangle$$

*(handwritten)* avoid by not letting A get too big:  $tr(A) \le C$

©Emily Fox 2014    33

5

# Mahalanobis Regret Minimization

*for Mahalanobis distance*

- Objective: $g_t^T A^{-1} g_t$

$$\min_A \sum_{t=1}^{T} \left\langle g_t, A^{-1} g_t \right\rangle \quad \text{subject to } A \succeq 0, \text{tr}(A) \leq C$$

- Solution:

$$A = c \left( \sum_{t=1}^{T} g_t g_t^T \right)^{\frac{1}{2}}$$

*if $Q$, $Q \succeq 0$, $\exists V$*
*$Q = V^T V$ ← square root matrix*

*Outer product of gradient*

For proof, see Appendix E, Lemma 15 of Duchi et al. 2011.
Uses "trace trick" and Lagrangian.

- *A* defines the norm of the metric space we should be operating in

34

---

# AdaGrad Algorithm

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w} \in \mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta A_t^{-1} g_t)||^2_{A_t}$$

- At time *t*, estimate optimal (sub)gradient modification *A* by

*estimate of A at time t*   *update time t*

$$A_t = \left( \sum_{\tau=1}^{t} g_\tau g_\tau^T \right)^{\frac{1}{2}}$$

*← in d dims*
*matrix $\sqrt{}$*
*is $O(d^3)$*

- For *d* large, $A_t$ is computationally intensive to compute. Instead,

$\text{diag}(A_t)$   $A_t = \begin{pmatrix} A_{ii} & 0 \\ 0 & \ddots \end{pmatrix}$

- Then, algorithm is a simple modification of normal updates:

$$\mathbf{w}^{(t+1)} = \arg\min_{\mathbf{w} \in \mathcal{W}} ||\mathbf{w} - (\mathbf{w}^{(t)} - \eta \, \text{diag}(A_t)^{-1} g_t)||^2_{\text{diag}(A_t)}$$

$$A_{ii}^t = \sqrt{\sum_{\tau=1}^{t} g_{\tau,i}^2}$$

*weigh dimensions by sqrt of sum of gradients in that dim*

35

6

# AdaGrad in Euclidean Space

$x_t = (0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ |v)$

- For $\mathcal{W} = \mathbb{R}^d$ ,

  $w^{(t+1)} \leftarrow w^{(t)} - \eta \, diag(A_t)^{-1} g_t$

  no constraints on w

- For each feature dimension,

  $$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta_{t,i} g_{t,i}$$

  adaptive step size

  where

  $$\eta_{t,i} = \eta / A_{t,ii}$$

  in sparse case, step size bigger when seeing a rare feature

- That is,

  $$w_i^{(t+1)} \leftarrow w_i^{(t)} - \frac{\eta}{\sqrt{\sum_{\tau=1}^{t} g_{\tau,i}^2}} g_{t,i}$$

- Each feature dimension has it's own learning rate!
  - Adapts with *t*
  - Takes geometry of the past observations into account
  - Primary role of η is determining rate the first time a feature is encountered

36

# AdaGrad Theoretical Guarantees

- AdaGrad regret bound:

  $$\sum_{t=1}^{T} f_t(\mathbf{w}^{(t)}) - f_t(\mathbf{w}^*) \leq 2R_\infty \sum_{i=1}^{d} ||g_{1:T,j}||_2$$

  $$R_\infty := \max_{t} ||\mathbf{w}^{(t)} - \mathbf{w}^*||_\infty$$

  radius of space

- So, what does this mean in practice?

- Many cool examples. This really is used in practice!
- Let's just examine one…

37

7

# AdaGrad Theoretical Example

$x = (0\ 000\ 1\ 0 - 0\ 10_l)$

- Expect to out-perform when gradient vectors are sparse
- SVM hinge loss example:

$$f_t(\mathbf{w}) = [1 - y^t \langle \mathbf{x}^t, \mathbf{w}\rangle]_+ \quad \text{where} \quad \mathbf{x}^t \in \{-1, 0, 1\}^d$$

hinge loss

- If $x_j^t \neq 0$ with probability $\propto j^{-\alpha}, \quad \alpha > 1$

$j^{-\alpha}$ → heavy tailed distribution

$R(T)/T$

$$\mathbb{E}\left[f\left(\frac{1}{T}\sum_{t=1}^{T}\mathbf{w}^{(t)}\right)\right] - f(\mathbf{w}^*) = \mathcal{O}\left(\frac{\|\mathbf{w}^*\|_\infty}{\sqrt{T}} \cdot \max\{\log d, d^{1-\alpha/2}\}\right)$$

for d small

- Previously best known method:

same

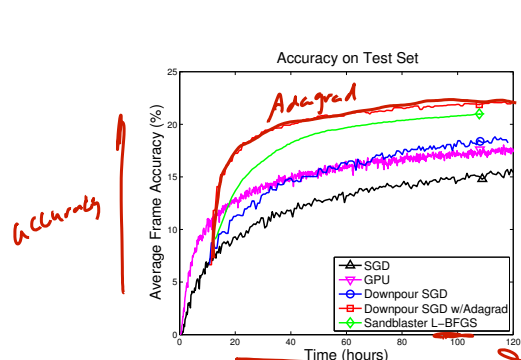$$\mathbb{E}\left[f\left(\frac{1}{T}\sum_{t=1}^{T}\mathbf{w}^{(t)}\right)\right] - f(\mathbf{w}^*) = \mathcal{O}\left(\frac{\|\mathbf{w}^*\|_\infty}{\sqrt{T}} \cdot \sqrt{d}\right)$$

adagrad can be exp better in d

©Emily Fox 2014

38

---

# Neural Network Learning

- Very non-convex problem, but use SGD methods anyway

$$f(\mathbf{w}, \xi) = \log\left(1 + \exp\left(\langle[p(\langle\mathbf{w}\right.\right.$$

$$p(\alpha) = \frac{1}{1 + \exp(\alpha)}$$

accuracy

**Accuracy on Test Set**

Adagrad

Average Frame Accuracy (%)

- △ SGD
- ▽ GPU
- ○ Downpour SGD
- □ Downpour SGD w/Adagrad
- ◇ Sandblaster L–BFGS

Time (hours)

running time

(Dean et al. 2012)

$p(\langle x_1, \xi_1\rangle)$

$w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5$

$\xi_1 \quad \xi_2 \quad \xi_3 \quad \xi_5 \quad \xi_4$

Distributed, $d = 1.7 \cdot 10^9$ parameters. SGD and AdaGrad use 80 machines (1000 cores), L-BFGS uses 800 (10000 cores)

*Images from Duchi et al. ISMP 2012 slides*

©Emily Fox 2014

39

8

# What you should know about Logistic Regression (LR) and Click Prediction

- Click prediction problem:
  - Estimate probability of clicking
  - Can be modeled as logistic regression
- Logistic regression model: Linear model
- Gradient ascent to optimize conditional likelihood
- Overfitting + regularization
- Regularized optimization
  - Convergence rates and stopping criterion
- Stochastic gradient ascent for large/streaming data
  - Convergence rates of SGD
- AdaGrad motivation, derivation, and algorithm

40