# Case Study 1: Estimating Click Probabilities

## Intro
## Logistic Regression
## Gradient Descent + SGD
## AdaGrad

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

January 7th, 2014

©Emily Fox 2014

1

---

# Ad Placement Strategies

- Companies bid on ad prices



- Which ad wins? (many simplifications here)
  - Naively:

  - But:

  - Instead:

©Emily Fox 2014

2

# Key Task: Estimating Click Probabilities

- What is the probability that user *i* will click on ad *j*

- Not important just for ads:
  - □ Optimize search results
  - □ Suggest news articles
  - □ Recommend products

- Methods much more general, useful for:
  - □ Classification
  - □ Regression
  - □ Density estimation
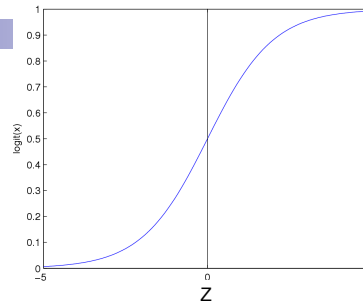
3

# Learning Problem for Click Prediction

- Prediction task:

- Features:

- Data:
  - □ Batch:

  - □ Online:

- Many approaches (e.g., logistic regression, SVMs, naïve Bayes, decision trees, boosting,…)
  - □ Focus on logistic regression; captures main concepts, ideas generalize to other approaches

4

# Logistic Regression

**Logistic function (or Sigmoid):** $\dfrac{1}{1 + exp(-z)}$

- Learn P(Y|**X**) directly
  - ☐ Assume a particular functional form
  - ☐ Sigmoid applied to a linear function of the data:

$$P(Y = 0|X, W) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

**Features can be discrete or continuous!**

5

---

# Very convenient!

$$P(Y = 0 \,|X = \, < X_1, ... X_n >) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

linear classification rule!

implies

$$\ln \frac{P(Y = 1 \,|X)}{P(Y = 0 \,|X)} = w_0 + \sum_i w_i X_i$$

6

# Digression: Logistic regression more generally

- Logistic regression in more general case, where *Y in* {y$_1$,...,y$_R$}

  for *k<R*
  $$P(Y = y_k|X) = \frac{\exp(w_{k0} + \sum_{i=1}^{n} w_{ki}X_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji}X_i)}$$

  for *k=R* (normalization, so no weights for this class)
  $$P(Y = y_R|X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^{n} w_{ji}X_i)}$$

  **Features can be discrete or continuous!**

7

# Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

- Discriminative (logistic regression) loss function:
  **Conditional Data Likelihood**

$$\ln P(\mathcal{D}_Y \mid \mathcal{D}_\mathbf{X}, \mathbf{w}) = \sum_{j=1}^{N} \ln P(y^j \mid \mathbf{x}^j, \mathbf{w})$$

8

# Expressing Conditional Log Likelihood

$$P(Y = 0|\mathbf{X}, \mathbf{w}) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$l(\mathbf{w}) \equiv \sum_j \ln P(y^j|\mathbf{x}^j, \mathbf{w})$$

$$P(Y = 1|\mathbf{X}, \mathbf{w}) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$\ell(\mathbf{w}) = \sum_j y^j \ln P(Y = 1|\mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(Y = 0|\mathbf{x}^j, \mathbf{w})$$

$$= \sum_j y^j (w_0 + \sum_{i=1}^d w_i x_i^j) - \ln \left( 1 + \exp(w_0 + \sum_{i=1}^d w_i x_i^j) \right)$$

9

# Maximizing Conditional Log Likelihood

$$l(\mathbf{w}) \equiv \ln \prod_j P(y^j|\mathbf{x}^j, \mathbf{w})$$

$$= \sum_j y^j (w_0 + \sum_{i=1}^d w_i x_i^j) - \ln \left( 1 + \exp(w_0 + \sum_{i=1}^d w_i x_i^j) \right)$$

Good news: $l(\mathbf{w})$ is concave function of $\mathbf{w}$, no local optima problems

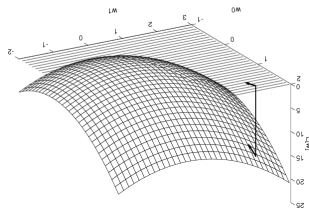Bad news: no closed-form solution to maximize $l(\mathbf{w})$

Good news: concave functions easy to optimize

10

5

# Optimizing concave function – Gradient ascent

- Conditional likelihood for Logistic Regression is concave. Find optimum with gradient ascent



**Gradient:** $\nabla_{\mathbf{w}} l(\mathbf{w}) = [\frac{\partial l(\mathbf{w})}{\partial w_0}, \ldots, \frac{\partial l(\mathbf{w})}{\partial w_n}]'$

Step size, $\eta > 0$

**Update rule:** $\triangle \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

- Gradient ascent is simplest of optimization approaches
  - □ e.g., Conjugate gradient ascent much better (see reading)

# Gradient Ascent for LR

Gradient ascent algorithm: iterate until change < ε

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$

For *i* = 1,…,*d*,

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$

repeat

# Regularized Conditional Log Likelihood

- If data is linearly separable, weights go to infinity
- Leads to overfitting → Penalize large weights

- Add regularization penalty, e.g., $L_2$:

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})) - \frac{\lambda}{2} ||\mathbf{w}||_2^2$$

- Practical note about $w_0$:

13

# Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg\max_{\mathbf{w}} \ln \left[ \prod_{j=1}^{N} P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})]$$

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg\max_{\mathbf{w}} \ln \left[ \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})) \right] - \frac{\lambda}{2} \sum_{i>0} w_i^2$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

14

# Stopping criterion

$$\ell(\mathbf{w}) = \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w})) - \frac{\lambda}{2} ||\mathbf{w}||_2^2$$

- Regularized logistic regression is strongly concave
  - Negative second derivative bounded away from zero:

- Strong concavity (convexity) is super helpful!!

- For example, for strongly concave *l*(**w**):

$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \frac{1}{2\lambda} ||\nabla \ell(\mathbf{w})||_2^2$$

15

# Convergence rates for gradient descent/ascent

- Number of Iterations to get to accuracy

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

- If func Lipschitz: $O(1/\epsilon^2)$

- If gradient of func Lipschitz: $O(1/\epsilon)$

- If func is strongly convex: $O(\ln(1/\epsilon))$

16

# Challenge 1: Complexity of computing gradients

- What's the cost of a gradient update step for LR???

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \sum_j x_i^j [y^j - \hat{P}(Y^j = 1 \mid \mathbf{x}^j, \mathbf{w}^{(t)})] \right\}$$

# Challenge 2: Data is streaming

- Assumption thus far: **Batch data**

- But, click prediction is a streaming data task:
  - □ User enters query, and ad must be selected:
    - Observe $\mathbf{x}^j$, and must predict $y^j$

  - □ User either clicks or doesn't click on ad:
    - Label $y^j$ is revealed afterwards
      - □ Google gets a reward if user clicks on ad

  - □ Weights must be updated for next time:

# Learning Problems as Expectations

- Minimizing loss in training data:
    - Given dataset:
        - Sampled iid from some distribution $p(\mathbf{x})$ on features:
    - Loss function, e.g., hinge loss, logistic loss,…
    - We often minimize loss in training data:

$$\ell_{\mathcal{D}}(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^{N} \ell(\mathbf{w}, \mathbf{x}^j)$$

- However, we should really minimize expected loss on all data:

$$\ell(\mathbf{w}) = E_{\mathbf{x}}\left[\ell(\mathbf{w}, \mathbf{x})\right] = \int p(\mathbf{x})\ell(\mathbf{w}, \mathbf{x})d\mathbf{x}$$

- So, we are approximating the integral by the average on the training data

19

# Gradient ascent in Terms of Expectations

- "True" objective function:

$$\ell(\mathbf{w}) = E_{\mathbf{x}}\left[\ell(\mathbf{w}, \mathbf{x})\right] = \int p(\mathbf{x})\ell(\mathbf{w}, \mathbf{x})d\mathbf{x}$$

- Taking the gradient:

- "True" gradient ascent rule:

- How do we estimate expected gradient?

20

## SGD: Stochastic Gradient Ascent (or Descent)

- "True" gradient:  $\nabla \ell(\mathbf{w}) = E_{\mathbf{x}}\left[\nabla \ell(\mathbf{w}, \mathbf{x})\right]$

- Sample based approximation:

- What if we estimate gradient with just one sample???
  - ☐ Unbiased estimate of gradient
  - ☐ Very noisy!
  - ☐ Called stochastic gradient ascent (or descent)
    - Among many other names
  - ☐ VERY useful in practice!!!

21

## Stochastic Gradient Ascent: general case

- Given a stochastic function of parameters:
  - ☐ Want to find maximum

- Start from $\mathbf{w}^{(0)}$
- Repeat until convergence:
  - ☐ Get a sample data point $\mathbf{x}^t$
  - ☐ Update parameters:

- Works on the online learning setting!
- Complexity of each gradient step is constant in number of examples!
- In general, step size changes with iterations

22

# Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_{\mathbf{x}}\left[\ell(\mathbf{w}, \mathbf{x})\right] = E_{\mathbf{x}}\left[\ln P(y|\mathbf{x}, \mathbf{w}) - \frac{\lambda}{2}||\mathbf{w}||_2^2\right]$$

- Batch gradient ascent updates:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left\{ -\lambda w_i^{(t)} + \frac{1}{N} \sum_{j=1}^{N} x_i^{(j)} [y^{(j)} - P(Y = 1|\mathbf{x}^{(j)}, \mathbf{w}^{(t)})] \right\}$$

- Stochastic gradient ascent updates:
  - Online setting:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta_t \left\{ -\lambda w_i^{(t)} + x_i^{(t)} [y^{(t)} - P(Y = 1|\mathbf{x}^{(t)}, \mathbf{w}^{(t)})] \right\}$$

---

# Convergence rate of SGD

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

  - Then, for step sizes:

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

# Convergence rates for gradient descent/ascent versus SGD

- Number of Iterations to get to accuracy

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

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

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

- 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

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

25

---

# Motivating AdaGrad (Duchi, Hazan, Singer 2011)

- Assuming $\mathbf{w} \in \mathbb{R}^d$, standard stochastic (sub)gradient descent updates are of the form:
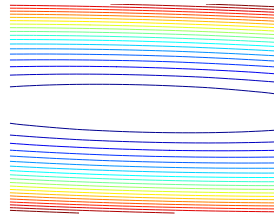
$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta g_{t,i}$$

- Should all features share the same learning rate?

- Often have high-dimensional feature spaces
  - Many features are irrelevant
  - Rare features are often very informative

- Adagrad provides a feature-specific adaptive learning rate by incorporating knowledge of the geometry of past observations
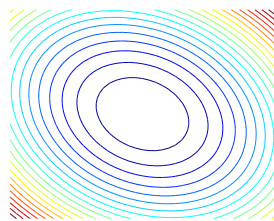
26

# Why Adapt to Geometry?

Hard

Nice

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

---

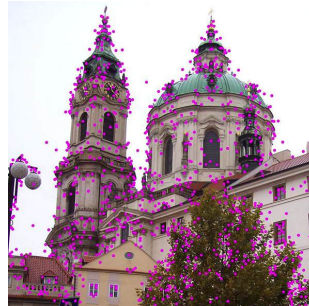# Not All Features are Created Equal

- 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.

High-dimensional image features

*Images from Duchi et al. ISMP 2012 slides*

©Emily Fox 2014

28

# Projected Gradient

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

- Brief aside…

- Consider an arbitrary feature space $\mathbf{w} \in \mathcal{W}$

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

    $$\mathbf{w}^{(t+1)} =$$

29

# Regret Minimization

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

- Algorithm iteratively predicts $\mathbf{w}^{(t)}$
- Incur ***loss*** $f_t(\mathbf{w}^{(t)})$
- ***Regret***:
  What is the total incurred loss of algorithm relative to the best choice of $\mathbf{w}$ that could have been made ***retrospectively***

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

30

# Regret Bounds for Standard SGD

- Standard projected gradient stochastic updates:

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

- Standard regret bound:

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

31

# 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

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

32

16

# Mahalanobis Regret Bounds

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

- **What *A* to choose?**

- Regret bound now:

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

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

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

©Emily Fox 2014     33

---

# Mahalanobis Regret Minimization

- Objective:

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

- Solution:

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

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

©Emily Fox 2014     34

17

# AdaGrad Algorithm

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

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

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

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

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

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

35

# AdaGrad in Euclidean Space

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

- For each feature dimension,

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

  where

$$\eta_{t,i} =$$

- 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

18

# AdaGrad Theoretical Guarantees

- AdaGrad regret bound:

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

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

- So, what does this mean in practice?

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

37

# AdaGrad Theoretical Example

- Expect to out-perform when gradient vectors are sparse
- SVM hinge loss example:
  $$f_t(\mathbf{w}) = [1 - y^t \langle \mathbf{x}^t, \mathbf{w} \rangle]_+ \quad \text{where} \quad \mathbf{x}^t \in \{-1, 0, 1\}^d$$

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

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

- Previously best known method:

$$\mathbb{E}\left[ f\left( \frac{1}{T} \sum_{t=1}^{T} \mathbf{w}^{(t)} \right) \right] - f(\mathbf{w}^*) = \mathcal{O}\left( \frac{||\mathbf{w}^*||_\infty}{\sqrt{T}} \cdot \sqrt{d} \right)$$
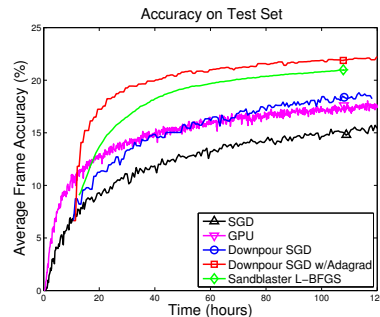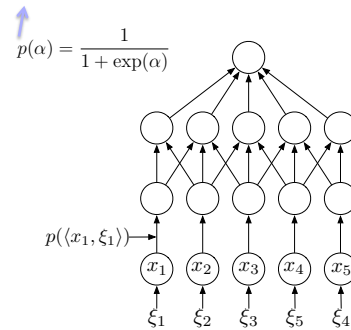
38

# Neural Network Learning

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

$$f(x;\xi) = \log\left(1 + \exp\left(\langle[p(\langle\text{...}\right.\right.$$

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

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

**Accuracy on Test Set**

(y-axis) Average Frame Accuracy (%)
(x-axis) Time (hours)

- SGD
- GPU
- Downpour SGD
- Downpour SGD w/Adagrad
- Sandblaster L–BFGS

(Dean et al. 2012)

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

---

# 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

©Emily Fox 2014

40