

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

$$C_1 \rightarrow \$10$$

$$C_2 \rightarrow \$20$$

$$C_3 \rightarrow \$100$$

- Which ad wins? (many simplifications here)

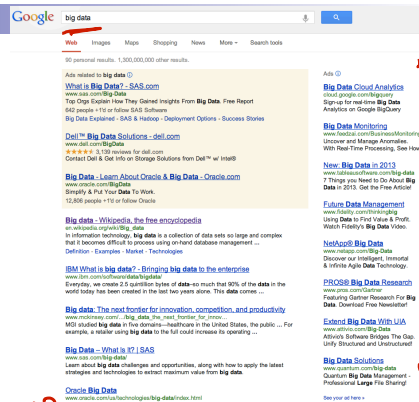
- Naively: $C_3 \rightarrow \$100$

- But: paid on clicks

- Instead:

$$e.g. P(\text{click} | C_3) = 0.01 \quad E[\$3] = 0.01 \times 100 = \$1$$

$$P(\text{click} | C_2) = 0.5 \quad E[\$2] = 0.5 \times 20 = \$10$$



©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

©Emily Fox 2014

3

Learning Problem for Click Prediction

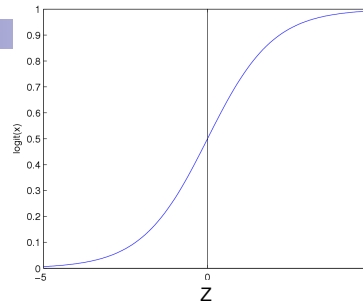
- Prediction task: $x \rightarrow [0,1]$ $p(\text{click}=1|x)$
- Features:
 $X = (\text{feats of page}, \text{feats ad}, \text{feats user})$
- Data: $(x, y) \rightarrow (\text{webpage 1}, \text{ad 7}, \text{user 25}, \text{time 12}) \rightarrow \text{click}=1$
 - Batch: fixed dataset $(x^1, y^1) \dots (x^N, y^N)$
 - Online: data as a stream
= user arrives at a page $\rightarrow x^t \rightarrow$ predict \hat{y}^t , click? \rightarrow observe y^t
- Many approaches (e.g., logistic regression, SVMs, naïve Bayes, decision trees, boosting,...)
 - Focus on logistic regression; captures main concepts, ideas generalize to other approaches

©Emily Fox 2014

4

Logistic Regression

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



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

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

linear function of features

Features can be discrete or continuous!

©Emily Fox 2014

5

Very convenient!

$$P(Y = 0 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

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

linear decision boundary

linear classification rule!

©Emily Fox 2014

6

Digression: Logistic regression more generally

- Logistic regression in more general case, where Y in $\{y_1, \dots, y_R\}$

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!

©Emily Fox 2014

7

Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

$$(x^i, y^i)_{i=1:N} \equiv \mathcal{D} = (\mathcal{D}_X, \mathcal{D}_Y)$$

- Discriminative (logistic regression) loss function:

Conditional Data Likelihood

$$\begin{aligned} \operatorname{argmax}_w P(\mathcal{D}_Y | \mathcal{D}_X, w) &\stackrel{\text{iid}}{\equiv} \operatorname{argmax}_w \prod_{j=1}^N P(y^j | x^j, w) \\ &\equiv \operatorname{argmax}_w \ln \prod_{j=1}^N P(y^j | x^j, w) \end{aligned}$$

$$\ln P(\mathcal{D}_Y | \mathcal{D}_X, w) = \sum_{j=1}^N \ln P(y^j | x^j, w)$$

goal to
max this

©Emily Fox 2014

8

Expressing Conditional Log Likelihood

argmax \mathbf{w}

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

$\{0,1\}$

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

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

$\equiv \sum_j \begin{cases} P(Y=1|\mathbf{x}^j, \mathbf{w}), & \text{if } y^j=1 \\ P(Y=0|\mathbf{x}^j, \mathbf{w}), & \text{if } y^j=0 \end{cases}$

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

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

argmax \mathbf{w} for LR

©Emily Fox 2014

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

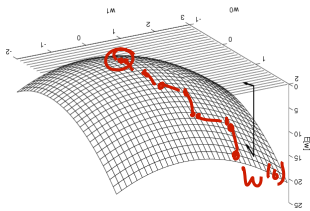
©Emily Fox 2014

10

Optimizing concave function — *f is concave*

Gradient ascent *diff. scale fnj convex f''(x) > 0 & x*

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



Gradient: $\nabla_{\mathbf{w}} l(\mathbf{w}) = \left[\frac{\partial l(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial l(\mathbf{w})}{\partial w_n} \right]^T$

Step size, $\eta > 0$

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

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

- 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 $< \epsilon$

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

truth (pointing to y^j) and *prediction* (pointing to \hat{P})

For $i = 1, \dots, d,$

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

weigh by feature value (pointing to x_i^j)

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 :

original w

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

conditional likelihood "regularization" penalty

$\sum_{i=1}^d w_i^2$

- Practical note about w_0 :

Don't regularize w_0 ...

©Emily Fox 2014

13

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

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

- Regularized maximum conditional likelihood estimate

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

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

neg derivative \leftarrow move towards zero

©Emily Fox 2014

14

Stopping criterion

for LR
strongly concave
with $\gamma \geq \lambda$

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

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

$f(x)$
concave $\Rightarrow -f''(x) \geq 0$; strongly concave: $-f''(x) \geq \gamma$
diff $\gamma > 0$

- Strong concavity (convexity) is super helpful!!

- For example, for strongly concave $\ell(\mathbf{w})$: stop $\ell(\mathbf{w}^*) - \ell(\mathbf{w}^t) \leq \epsilon$

$$\underbrace{\ell(\mathbf{w}^*)}_{\text{opt}} - \underbrace{\ell(\mathbf{w})}_{\text{value at some } t} \leq \frac{1}{2\lambda} \|\nabla \ell(\mathbf{w})\|_2^2 \leq \epsilon \Rightarrow \|\nabla \ell(\mathbf{w}^t)\|_2^2 \leq 2\lambda \epsilon$$

©Emily Fox 2014

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: $O(1/\epsilon^2)$
 $\|\ell(u) - \ell(v)\| \leq K \|u - v\|$
- If gradient of func Lipschitz: $O(1/\epsilon)$
 $\|\nabla \ell(u) - \nabla \ell(v)\| \leq K \|u - v\|$
- If func is strongly convex: $O(\ln(1/\epsilon))$

regularized LR

exp
faster

(sufficiently large)
constant
step size

©Emily Fox 2014

16

Challenge 1: Complexity of computing gradients

d features

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

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

for each i

O(Nd)

forall features i

cost is O(Nd^2)

cache P(Y^j=1|x^j,w^t) O(Nd)

in Big Data N is very large:

O(Nd) → only taking little η step

©Emily Fox 2014

17

Challenge 2: Data is streaming

- Assumption thus far: Batch data

$$\sum_{j=1}^N \dots$$

- But, click prediction is a streaming data task:

- User enters query, and ad must be selected:

- Observe x^i , and must predict y^i

q → [] → x^i → predict g^i → show ad

- User either clicks or doesn't click on ad:

- Label y^i is revealed afterwards
- Google gets a reward if user clicks on ad

- Weights must be updated for next time:

$$w^{(t+1)} \leftarrow w^{(t)} + \Delta$$

©Emily Fox 2014

18

Learning Problems as Expectations

- Minimizing loss in training data:

- Given dataset: $x^1 \dots x^N$
 - Sampled iid from some distribution $p(\mathbf{x})$ on features: $x^i \sim p(\mathbf{x})$
- Loss function, e.g., hinge loss, logistic loss,...
- We often minimize loss in training data:

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

Monte Carlo integration

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

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

true expected loss

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

Gradient ascent in Terms of Expectations

- “True” objective function:

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

- Taking the gradient:

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

- “True” gradient ascent rule:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta E_{\mathbf{x}} [\nabla_{\mathbf{w}} \ell(\mathbf{w}, \mathbf{x})]$$

- How do we estimate expected gradient? *estimate exp. gradient from data*

if data on machine, must shuffle first

SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient: $\nabla l(\mathbf{w}) = E_{\mathbf{x}} [\nabla l(\mathbf{w}, \mathbf{x})]$

- Sample based approximation: $\mathbf{x} \stackrel{iid}{\sim} P(\mathcal{X})$

$$\nabla l(\mathbf{w}) = E_{\mathbf{x}} [\nabla l(\mathbf{w}, \mathbf{x})] \approx \hat{\nabla} l(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \nabla l(\mathbf{w}, \mathbf{x}^j)$$

the bigger N , the closer $\hat{\nabla} l$ is ∇l

- What if we estimate gradient with just one sample???

- Unbiased estimate of gradient $\nabla l(\mathbf{w}) \approx \hat{\nabla} l(\mathbf{w}) = \nabla l(\mathbf{w}, \mathbf{x}^{(t)})$
- Very noisy! $E[\hat{\nabla} l(\mathbf{w})] = E_{\mathbf{x}^{(t)}}[\nabla l(\mathbf{w}, \mathbf{x}^{(t)})] = \nabla l(\mathbf{w})$
- Called stochastic gradient ascent (or descent)
 - Among many other names
- VERY useful in practice!!!

©Emily Fox 2014

21

Stochastic Gradient Ascent: general case

running avg of $w^{(t)}$ is used for prediction

- Given a stochastic function of parameters: $f(\mathbf{w}) = E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$

- Want to find maximum

$$\mathbf{w}^* \in \underset{\mathbf{w}}{\operatorname{argmin}} E_{\mathbf{x}} [f(\mathbf{w}, \mathbf{x})]$$

- Start from $\mathbf{w}^{(0)}$ e.g., $\mathbf{w}^{(0)} = \mathbf{0}$

- Repeat until convergence:

- Get a sample data point \mathbf{x}^t
 Sell and, assume $y_k \in \text{click?}$
- Update parameters:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta_t \nabla f(\mathbf{w}^{(t)}, \mathbf{x}^t)$$

actual params, not avg

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

$$\text{e.g. } \eta_t = \frac{K}{t} \text{ for } K > 0$$

Stopping criteria?

good luck...

in online: never stop

as batch alg:

- validation set

- progressive validation

©Emily Fox 2014

22

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

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

- Batch gradient ascent updates: $O(nd)$

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

- Stochastic gradient ascent updates:

- Online setting:

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

1 data point at a time