

Stochastic Gradient Descent

Machine Learning – CSE546
Sham Kakade
University of Washington
October 18, 2016
©Sham Kakade 2016

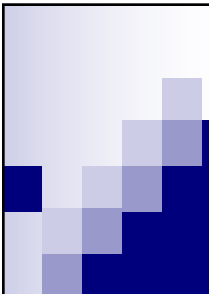
1

Announcements:

- HW2 posted today.
- Today:
 - Review: Logistic Regression
 - New: GD, SGD

©2016 Sham Kakade


2



Classification Logistic Regression

Machine Learning – CSE546
Sham Kakade
University of Washington
October 13, 2016
©Sham Kakade 2016

3



**THUS FAR, REGRESSION:
PREDICT A CONTINUOUS
VALUE GIVEN SOME INPUTS**

©Sham Kakade 2016

4

Weather prediction revisited

©Sham Kakade 2016 5

Reading Your Brain, Simple Example

[Mitchell et al.]
Pairwise classification accuracy: 85%

©Sham Kakade 2016 6

Classification

X - image
Y = {flower, not flower}
Y = {0, 1} Binary class

- Learn: $h: \mathbf{X} \mapsto Y$
 - \mathbf{X} - features
 - Y - target classes
- Conditional probability: $P(Y|\mathbf{X})$
- Suppose you know $P(Y|\mathbf{X})$ exactly, how should you classify?
 - Bayes optimal classifier:
 $\hat{y} = \arg \max_y P_r(Y=y|\mathbf{X})$
- How do we estimate $P(Y|\mathbf{X})$? *for '0/1 loss'*

©Sham Kakade 2016 7

Link Functions

- Estimating $P(Y|\mathbf{X})$: Why not use standard linear regression?
 $Y \approx w_0 + \sum w_i X_i$
- Combining regression and probability?
 - Need a mapping from real values to $[0, 1]$
 - A link function!

©Sham Kakade 2016 8

Logistic Regression

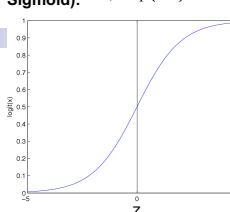
Logistic function (or Sigmoid): $\frac{1}{1 + \exp(-z)}$

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

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

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

Features can be discrete or continuous!

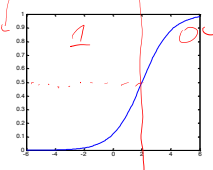


Understanding the sigmoid

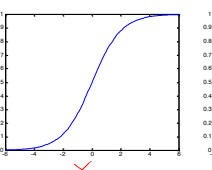
$P(Y=0|X,w) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$

$P(Y=1|X,w)$

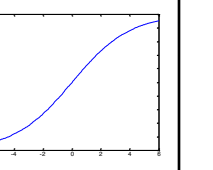
$w_0 = -2, w_1 = -1$



$w_0 = 0, w_1 = -1$



$w_0 = 0, w_1 = -0.5$

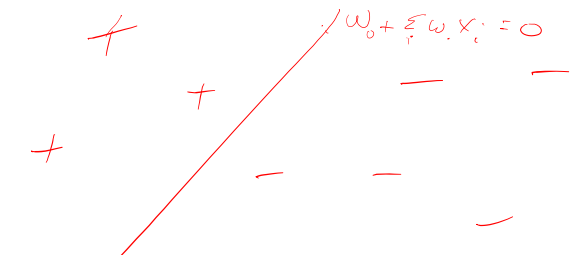


Logistic Regression – a Linear classifier

$\frac{1}{1 + \exp(-z)}$

$$g(w_0 + \sum_i w_i x_i) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

$w_0 + \sum_i w_i x_i = 0$



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

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

implies

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

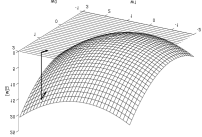
implies

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

linear classification rule!

Optimizing concave function – Gradient ascent

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



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

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 can be much better

©Sham Kakade 2016

13

Loss function: Conditional Likelihood

- Have a bunch of iid data of the form:

assume $(x^1, y^1), \dots, (x^N, y^N)$
 $i=1 \text{ to } N$
 $\mathbb{R}^n \times \mathbb{R}$

- Discriminative (logistic regression) loss function:

Conditional Data Likelihood

$$\arg \max_{\mathbf{w}} P(y^1, \dots, y^N | x^1, \dots, x^N, \mathbf{w})$$

$$= \arg \max_{\mathbf{w}} \prod_j P(y^j | x^j, \mathbf{w})$$

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

©Sham Kakade 2016

14

Expressing Conditional Log Likelihood

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

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

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

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

$$= \sum_j y^j \ln \frac{e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}} + (1 - y^j) \ln \frac{1}{1 + e^{w_0 + \sum_i w_i x_i^j}}$$

$$= \sum_{j=1}^N y^j (w_0 + \sum_i w_i x_i^j) - \ln (1 + e^{w_0 + \sum_i w_i x_i^j})$$

©Sham Kakade 2016

15

Maximizing Conditional Log Likelihood

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

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

$$l(\mathbf{w}) \equiv \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w})$$

$$= \sum_{j=1}^N y^j (w_0 + \sum_i w_i x_i^j) - \ln (1 + \exp(w_0 + \sum_i w_i x_i^j))$$

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

©Sham Kakade 2016

16

Maximize Conditional Log Likelihood: Gradient ascent

$$l(\mathbf{w}) = \sum_j y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))$$

$$\frac{\partial l}{\partial w_k} = \sum_j y^j x_k^j - \frac{x_k^j \exp(\eta)}{1 + \exp(\eta)}$$

$$= \sum_j x_k^j (y^j - \hat{P}(y^j=1 | x^j, \mathbf{w}))$$

$$\frac{\partial l}{\partial w_0} = \sum_j (y^j - \hat{P}(y^j=1 | x^j, \mathbf{w}))$$

©Sham Kakade 2016

17

Gradient Ascent for LR

Gradient ascent algorithm: iterate until change $< \epsilon$

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

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

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

repeat

©Sham Kakade 2016

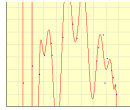
18

Regularization in linear regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$

$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



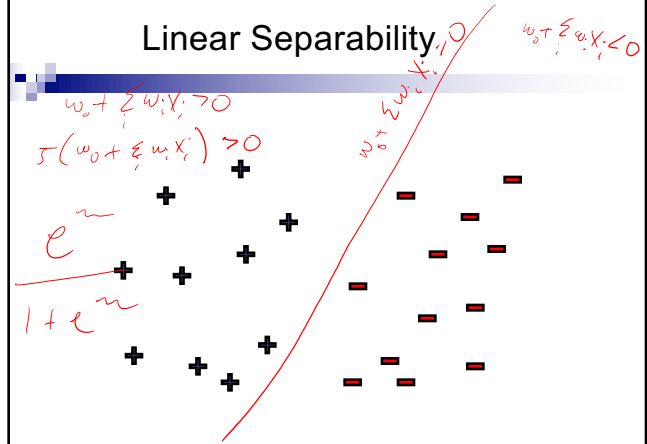
- Regularized least-squares (a.k.a. ridge regression), for $\lambda > 0$:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j (t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j))^2 + \lambda \sum_{i=1}^k w_i^2$$

©Sham Kakade 2016

19

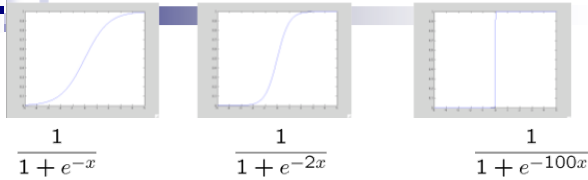
Linear Separability



©Sham Kakade 2016

20

Large parameters → Overfitting



- If data is linearly separable, weights go to infinity
- In general, leads to overfitting:
- Penalizing high weights can prevent overfitting...

©Sham Kakade 2016

21

Regularized Conditional Log Likelihood

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

$$\ell(\mathbf{w}) = \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Practical note about w_0 :

don't regularize w_0

- Gradient of regularized likelihood:

©Sham Kakade 2016

22

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

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

- Regularized maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) - \frac{\lambda}{2} \sum_{i=1}^k w_i^2$$

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

©Sham Kakade 2016

23

Please Stop!! Stopping criterion

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

- When do we stop doing gradient descent?

- Because $\ell(\mathbf{w})$ is strongly concave:

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

- Thus, stop when:

gradient is small

©Sham Kakade 2016

24

Convergence rates for gradient descent/ascent

- Number of Iterations to get to accuracy

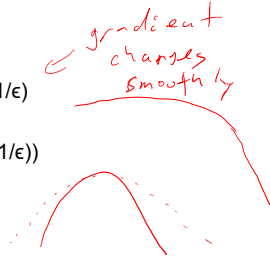
$$\ell(\mathbf{w}^*) - \ell(\mathbf{w}) \leq \epsilon$$

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

② If gradient of func Lipschitz: $O(1/\epsilon)$

③ If func is strongly convex: $O(\ln(1/\epsilon))$

② or ③ constant learning rate



©Sham Kakade 2016

25

Digression: Logistic regression for more than 2 classes

- Logistic regression in more general case (C classes), where $Y \in \{0, \dots, C-1\}$

©Sham Kakade 2016

26

Digression: Logistic regression more generally

- Logistic regression in more general case, where $Y \in \{0, \dots, C-1\}$

for $c > 0$

$$P(Y = c | \mathbf{x}, \mathbf{w}) = \frac{\exp(w_{c0} + \sum_{i=1}^k w_{ci}x_i)}{1 + \sum_{c'=1}^{C-1} \exp(w_{c'0} + \sum_{i=1}^k w_{c'i}x_i)}$$

for $c=0$ (normalization, so no weights for this class)

$$P(Y = 0 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \sum_{c'=1}^{C-1} \exp(w_{c'0} + \sum_{i=1}^k w_{c'i}x_i)}$$

Learning procedure is basically the same as what we derived!

©Sham Kakade 2016

27

Stochastic Gradient Descent

Machine Learning – CSE546

Sham Kakade

University of Washington

October 18, 2016

©Sham Kakade 2016

28

The Cost, The Cost!!! Think about the cost... $\sum (y - \hat{y})$ $x^j \in \mathbb{R}^d$

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

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

cost of updating 4 - coordinate $O(Nd)$
 Naively $O(Nd^2)$ for all-coordinates
 Re-use \hat{P} computation $O(Nd)$

©Sham Kakade 2016

29

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:

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

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

©Sham Kakade 2016

30

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:
- "True" gradient ascent rule:
- How do we estimate expected gradient?

©Sham Kakade 2016

31

SGD: Stochastic Gradient Ascent (or Descent)

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

©Sham Kakade 2016

32

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

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

- 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:

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

©Sham Kakade 2016

33

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

©Sham Kakade 2016

34

What you should know...

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
 - Logistic function maps real values to $[0, 1]$
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Cost of gradient step is high, use stochastic gradient descent

©Sham Kakade 2016

35