

S&DS 355 / 365 / 565
Data Mining and Machine Learning

Gradient Descent

Thursday, September 12

Yale

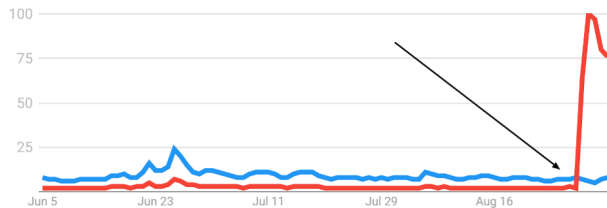
Outline

- Administrivia
- Recap
- Shrinkage, bias and variance
- Regularization
- Stochastic gradient descent

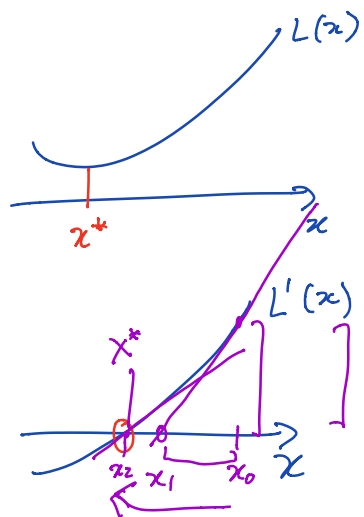
Administrivia

- Classes will split either next week or the week after
- 365 Quiz postponed

ML in the news



Google searches for "BTC" (red) vs "bitcoin" (blue) searches

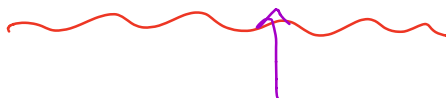


$$L'(x^*) = 0$$

$$L''(x_0) = \frac{L'(x_0)}{x_0 - x_1}$$

$$x_0 - x_1 = \frac{L'(x_0)}{L''(x_0)}$$

$$x_1 = x_0 - \left(\frac{1}{L''(x_0)} \right) L'(x_0)$$



Fitting a logistic regression model

- We maximize conditional likelihood. There is no closed form.
- Need to iterate.
- Standard approach is equivalent to Newton's algorithm
 - ▶ Make a quadratic approximation
 - ▶ Do a weighted least squares regression
 - ▶ Repeat

Newton's method

To find a zero of $f(x)$:

$$x \longleftarrow x - \frac{f(x)}{f'(x)}$$

To find a maximum of $f(x)$:

$$x \longleftarrow x - H(f, x)^{-1} \nabla f(x)$$

where ∇f is the (gradient) vector of first, derivatives, and H is the (Hessian) matrix of second derivatives

Iteratively reweighted least squares

(not in syllabus)

Given the current estimate $\hat{\beta}$, Newton's algorithm forms a quadratic approximation to the log-likelihood:

$$-\ell(\beta) = \frac{1}{2}(z - X\beta)^T W(z - X\beta) + \text{constant}.$$

where

$$z_i = \log \left(\frac{\pi_1(x_i)}{1 - \pi_1(x_i)} \right) + \frac{y_i - \pi_1(x_i)}{\pi_1(x_i)(1 - \pi_1(x_i))}.$$

is a “synthetic” response.

W is a diagonal weight matrix, with weight on the i th point given by

$$w_i = \pi_1(x_i)(1 - \pi_1(x_i))$$

This is a weighted least squares problem.

Discriminative vs	Generative
Linear Regression	LDA / QDA
Logistic Regression	
kNN	

Penalization, shrinkage, bias and variance

Estimator $\hat{\theta}$ of a parameter θ :

$$\begin{array}{ll} \text{bias}^2 & \left(\mathbb{E}(\hat{\theta}) - \theta \right)^2 \\ \text{variance} & \mathbb{E}(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2 \end{array}$$

Penalization, shrinkage, bias and variance

Estimator $\hat{\theta}$ of a parameter θ :

$$\begin{array}{ll} \text{bias}^2 & \left(\mathbb{E}(\hat{\theta}) - \theta\right)^2 \\ \text{variance} & \mathbb{E}(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2 \end{array}$$

Expected squared error decomposes as

$$\mathbb{E}(\hat{\theta} - \theta)^2 = \text{bias}^2 + \text{variance}$$

Penalization, bias and variance

Let's see how shrinkage affects the bias and variance.

Suppose $Y \sim N(\theta, \sigma^2)$.

(a) $\hat{\theta} = Y$. Bias? Variance?

$$E(\hat{\theta}) - \theta = 0 \quad \text{Var}(\hat{\theta}) = \sigma^2$$

Penalization, bias and variance

Let's see how shrinkage affects the bias and variance.

Suppose $Y \sim N(\theta, \sigma^2)$.

(a) $\hat{\theta} = Y$. Bias? Variance?

(b) $\hat{\theta} = bY$, for $0 \leq b \leq 1$. Bias? Variance?

$$\mathbb{E}(\hat{\theta}) - \theta = b\theta - \theta = ((b-1)\theta)^2$$

$$\text{Var}(\hat{\theta}) = b^2 \sigma^2$$

Exercise

Consider the simplified version of the objective function

$$F(\beta) = (Y - \beta)^2 + \lambda\beta^2$$

What is the minimizer $\hat{\beta}$?

Exercise

$$Y \sim N(\beta, \sigma^2)$$

Consider the simplified version of the objective function

$$\rightarrow F(\beta) = \underbrace{(Y - \beta)^2}_{\text{wavy}} + \underbrace{\lambda \beta^2}_{\text{boxed and crossed out}}$$

What is the minimizer $\hat{\beta}$?

$$\boxed{\hat{\beta} = \left(\frac{1}{1 + \lambda} \right) Y}$$

(bY)

Exercise

Now let's add a predictor variable,

$$F(\beta) = (Y - X\beta)^2 + \lambda\beta^2$$

What is the minimizer?

Exercise

Now let's add a predictor variable,

$$F(\beta) = (Y - X\beta)^2 + \lambda\beta^2$$

What is the minimizer?

$$\hat{\beta} = \frac{XY}{X^2 + \lambda}$$

Shrinkage

In a *shrinkage estimator*, we squash down the estimate by a scaling factor. For example,

$$\hat{\beta} \leftarrow \left(\frac{1}{1 + \lambda} \right) \hat{\beta}$$

This induces bias-variance tradeoff — the bias goes up, but the variance goes down.

Penalization

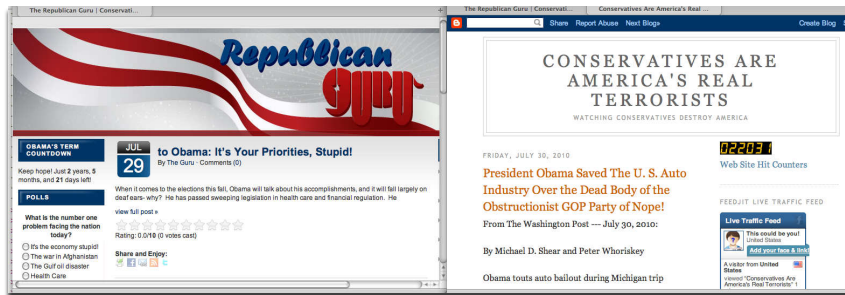
To guard against overfitting, we can *penalize* the coefficients:

$$F(\beta) = -\log\text{-likelihood}(\beta) + \lambda \|\beta\|^2$$

- Large coefficients incur a large penalty
- The *regularization parameter* λ controls the tradeoff between fit to the data, and size of the coefficients
- Small λ : high variance, low bias
- Large λ : low variance, high bias
- As λ increases, the size of the coefficients $|\beta_j|$ decreases.

Political blog classification

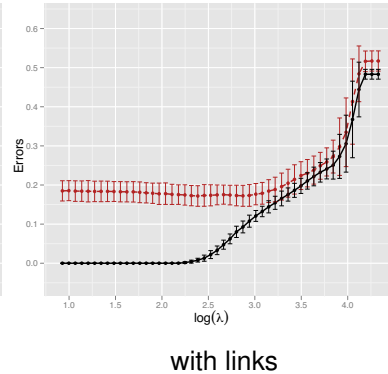
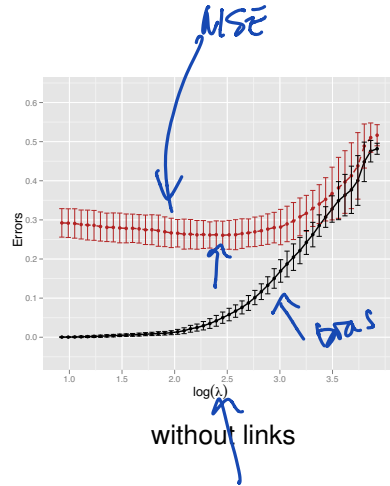
- Political Blog Classification. A collection of 403 political blogs were collected during two months before the 2004 presidential election. The goal is to predict whether a blog is *liberal* ($Y = 0$) or *conservative* ($Y = 1$) given the content of the blog.



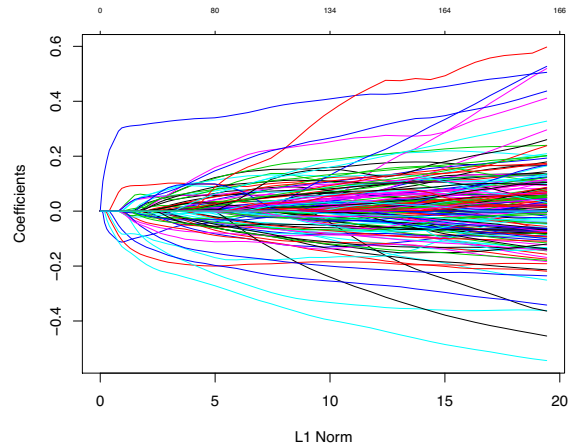
Political blog classification

- 403 blogs
- 205 are “liberal” and 198 are “conservative”
- For each word, value of a feature is word frequency
- Lower case and remove highly frequent words, throw out those appearing fewer than 10 times.
- 23,955 features
- Links to 292 popular blogs included as binary vector

Political blog classification results

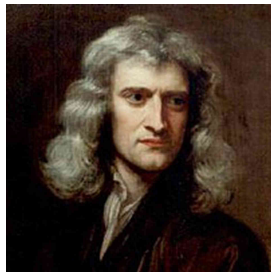


Political blog classification results



regularization paths

Newton's method

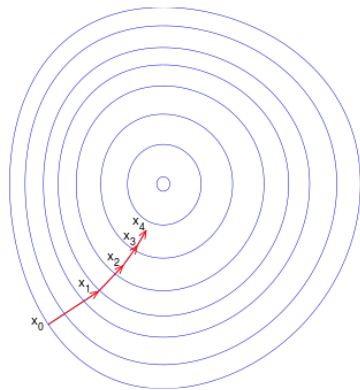


- “Grand-daddy of optimization”
- Fast convergence
- Great properties
- Second order method
- Not scalable

Gradient Descent

Gradient descent is a procedure for finding the arguments that minimize a particular function (called a **cost function**).

e.g. cost function could be a negative likelihood or negative log-likelihood.



Gradient Descent

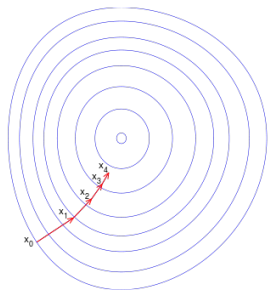
Goal: Find $(\theta_1, \dots, \theta_p)$ that minimizes **cost function** $L(\theta_1, \dots, \theta_p)$

Update equation:

$$\theta_j \leftarrow \theta_j - \rho \frac{\partial L}{\partial \theta_j}$$

In matrix form:

$$\theta \leftarrow \theta - \rho \nabla L(\theta)$$



ρ is called the learning rate.

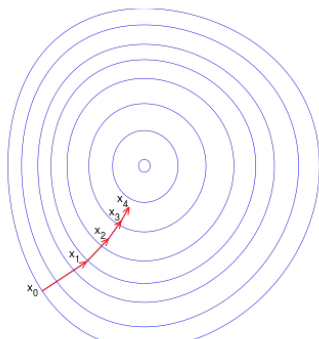


Intuition for GD

The gradient gives the direction of *steepest ascent*. Consider the Taylor Series of f about θ ,

$$f(\theta + \delta) = f(\theta) + \nabla f(\theta) \cdot \delta + \dots,$$

for some small increment δ . Then, ignoring higher order terms, it is clear that to maximize $f(\theta + \delta)$ we must pick δ in the direction of $\nabla f(\theta)$.



GD for Logistic Regression

Cost function (negative log-likelihood):

$$L(\beta) = - \sum \left(y_i x_i^t \beta - \log(1 + e^{x_i^t \beta}) \right)$$

Update for logistic regression:

$$\beta_j \leftarrow \beta_j - \rho \frac{\partial L}{\partial \beta_j}.$$

$$\beta_j \leftarrow \beta_j - \rho \sum_i \left(\frac{e^{x_i^t \beta}}{1 + e^{x_i^t \beta}} - y_i \right) x_{ij}.$$

GD for Linear Regression

Gradient descent also works for linear regression:

$$L(\beta) = \|Y - X\beta\|^2$$

$$\frac{\partial L(\beta)}{\partial \beta_j} = -2 \sum_{i=1}^n (y_i - x_i^t \beta) x_{ij}$$

GD update step:

$$\beta_j \leftarrow \beta_j + \rho \sum_{i=1}^n (y_i - x_i^t \beta) x_{ij}.$$

GD for Linear Regression

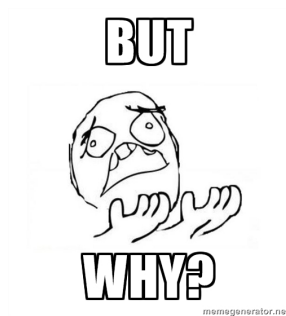
Gradient descent also works for linear regression:

$$L(\beta) = \|Y - X\beta\|^2$$

$$\frac{\partial L(\beta)}{\partial \beta_j} = -2 \sum_{i=1}^n (y_i - x_i^t \beta) x_{ij}$$

GD update step:

$$\beta_j \leftarrow \beta_j + \rho \sum_{i=1}^n (y_i - x_i^t \beta) x_{ij}.$$



Improving Upon Gradient Descent

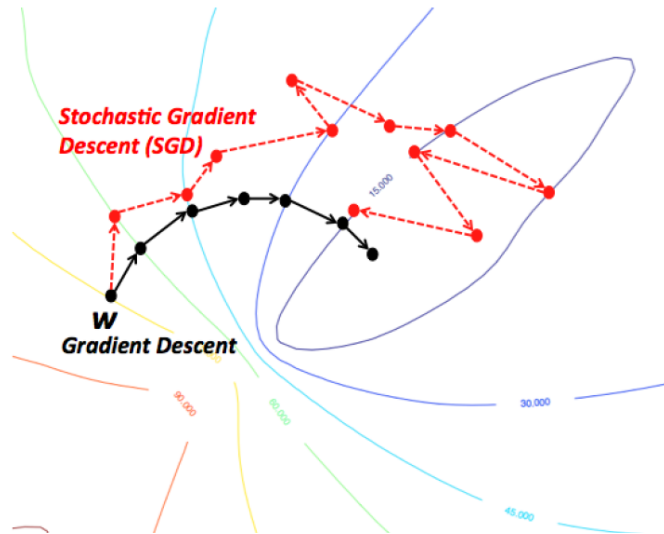
Each step of (batch) gradient descent requires a calculation involving all of the data points.

Improving Upon Gradient Descent

Each step of (batch) gradient descent requires a calculation involving all of the data points.

Stochastic gradient descent, in contrast, only computes based on a smaller subset of the data points (e.g. 1 observation) at each step.

Diagrammatic Differences



II: Stochastic gradient descent

- Suppose that we want to fit a really big model, n and d very large
- What is complexity of least squares and IRLS?

$$\underbrace{(X^T X)^{-1} X^T y}$$

$$\begin{matrix} X^T X \\ (d \times d) \end{matrix} \quad O(d^3) \quad + O(n d^2) \quad \begin{matrix} n \times d \times d \end{matrix}$$

II: Stochastic gradient descent

- Suppose that we want to fit a really big model, n and d very large
- What is complexity of least squares and IRLS?

$$O(nd^2 + d^3)$$

- Not practical for large problems
- Second order method, using covariance matrix and Hessian, solving linear system
- How can we get this to scale?

Typical example

- We want to classify documents according to whether or not they are about corporate news.
- There are about 780,000 documents in the collection
- 60,000,000 words
- Document represented as sparse tf-idf vector

1 | 5 : 1.1789641e-01 39 : 6.0373064e-02 45 : 1.3163488e-01

- The text takes about 1.1 Gbytes
- How can we efficiently train a classifier?

Online learning

We will introduce a method that

- Reads in the documents one at a time
- Updates the model for each document
- Updates are linear in the size of the document
- Uses little memory, never reads in the entire corpus
- Only stores a dictionary of feature weights

Stochastic gradient descent

We initialize all weights to zero: $\beta_j = 0, j = 1, \dots, d$.

We read through the data one record at a time, and update the model.

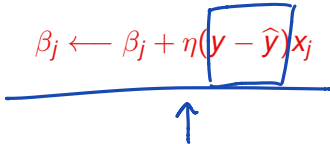
- 1 Read data item x
- 2 Make a prediction $\hat{y}(x) = \sum_{j=1}^d \beta_j x_j$
- 3 Observe the true response/label y
- 4 Update the weights β so \hat{y} is closer to y

Stochastic gradient descent

To begin, suppose we are just doing *linear regression*. We initialize all weights to zero: $\beta_j = 0, j = 1, \dots, d$.

We read through the data one record at a time, and update the model.

- 1 Read data item x
- 2 Make a prediction $\hat{y}(x) = \sum_{j=1}^d \beta_j x_j$
- 3 Observe the true response/label y
- 4 Update the weights β so \hat{y} is closer to y



The diagram shows the weight update equation $\beta_j \leftarrow \beta_j + \eta(y - \hat{y})x_j$ written in red. A blue box highlights the term $(y - \hat{y})x_j$. A blue arrow points upwards from below the line to the box. A blue bracket is positioned to the right of the list of steps, spanning from step 2 to step 4.

$$\beta_j \leftarrow \beta_j + \eta(y - \hat{y})x_j$$

Stochastic gradient descent

Think about how to apply this to the data we discussed earlier

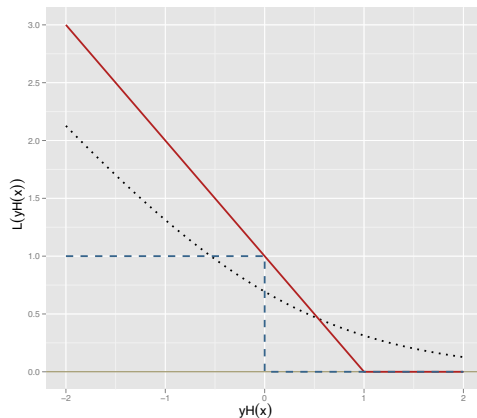
- We want to classify documents according to whether or not they are about corporate news.
- There are about 780,000 documents in the collection
- 60,000,000 words
- Document represented as sparse tf-idf vector

1 | 5 : 1.1789641e-01 39 : 6.0373064e-02 45 : 1.3163488e-01

- The text takes about 1.1 Gbytes

Key is that each data item is represented sparsely

SGD for general loss




- Squared error $L(y, \beta^T x) = (y - x^T \beta)^2$
- Logistic: $L(y, \beta^T x) = -y x^T \beta + \log(1 + \exp(x^T \beta))$
- Hinge: $L(y, \beta^T x) = (1 - y x^T \beta)_+$

SGD for general loss

SGD update:

$$\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \eta \nabla L(\boldsymbol{y}, \boldsymbol{\beta}^T \boldsymbol{x})$$

$$\beta_j \leftarrow \beta_j - \eta \frac{\partial L(\boldsymbol{y}, \boldsymbol{\beta}^T \boldsymbol{x})}{\partial \beta_j}$$

- 
- η is the *learning rate* or “step size”
 - Needs to be chosen carefully, getting smaller over time

SGD for general loss

SGD update:

$$\boldsymbol{\beta} \longleftarrow \boldsymbol{\beta} - \eta \nabla L(\boldsymbol{y}, \boldsymbol{\beta}^T \boldsymbol{x})$$

$$\beta_j \longleftarrow \beta_j - \eta \frac{\partial L(\boldsymbol{y}, \boldsymbol{\beta}^T \boldsymbol{x})}{\partial \beta_j}$$

Some intuition for what this is doing, and why the step size needs to decrease

$$L(\boldsymbol{\beta} + \epsilon \boldsymbol{v}) \approx L(\boldsymbol{\beta}) + \epsilon \boldsymbol{v}^T \nabla L(\boldsymbol{\beta})$$

$$L(\boldsymbol{\beta} - \eta \nabla L(\boldsymbol{\beta})) \approx L(\boldsymbol{\beta}) - \eta \|\nabla L(\boldsymbol{\beta})\|^2$$

This is why SGD is going downhill — if η is small enough

SGD for logistic regression

SGD Update:

$$\beta_j \leftarrow \beta_j + \eta(y - \pi)x_j$$

]

$$\beta_j x_j \leftarrow \beta_j x_j + \eta(\overset{\downarrow}{y} - \pi) \overset{0.7}{x_j^2}$$

$$\pi = \frac{1}{1 + \exp(-\beta^T \mathbf{x})}$$

]

Case checking:

$$\mathbb{P}(Y=1 | X=x)$$

- Suppose $y = 1$ and probability π is high?

SGD for logistic regression

SGD Update:

$$\beta_j \leftarrow \beta_j + \eta(y - \pi)x_j$$

$$\uparrow \quad \textcircled{\beta_j x_j} \leftarrow \beta_j x_j + \overset{1}{\eta} \overset{0.1}{(y - \pi)} x_j^2$$

$$\pi = \frac{1}{1 + \exp(-\underline{\beta^T \mathbf{x}})}$$

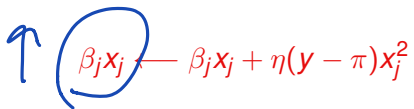
Case checking:

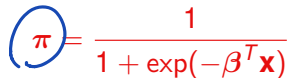
- Suppose $y = 1$ and probability π is high? *small change*
- Suppose $y = 1$ and probability π is small?

SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - \pi)x_j$$


$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - \pi)x_j^2$$


$$\pi = \frac{1}{1 + \exp(-\beta^T \mathbf{x})}$$

Case checking:

- Suppose $y = 1$ and probability π is high? *small change*
- Suppose $y = 1$ and probability π is small? *big change* \uparrow
- Suppose $y = 0$ and probability π is small?

SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - \pi)x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - \pi)x_j^2$$

$$\pi = \frac{1}{1 + \exp(-\beta^T \mathbf{x})}$$

Case checking:

- Suppose $y = 1$ and probability π is high? *small change*
- Suppose $y = 1$ and probability π is small? *big change* \uparrow
- Suppose $y = 0$ and probability π is small? *small change*
- Suppose $y = 0$ and probability π is big?

SGD for logistic regression

SGD Update:

$$\beta_j \longleftarrow \beta_j + \eta(y - \pi)x_j$$

$$\beta_j x_j \longleftarrow \beta_j x_j + \eta(y - \pi)x_j^2$$

$$\pi = \frac{1}{1 + \exp(-\beta^T \mathbf{x})}$$

Case checking:

- Suppose $y = 1$ and probability π is high? *small change*
- Suppose $y = 1$ and probability π is small? *big change* \uparrow
- Suppose $y = 0$ and probability π is small? *small change*
- Suppose $y = 0$ and probability π is big? *big change* \downarrow

SGD: choice of learning rate

A conservative choice of learning rate is

$$\eta_t = \frac{1}{t}$$

A more aggressive choice is

$$\eta_t = \frac{1}{\sqrt{t}}$$

Which is more appropriate for GD?

Which is more appropriate for SGD?

SGD: choice of learning rate

Learning rate should scale as

$$\eta_t = \frac{1}{\sqrt{t}}$$

Problem: Some of the updates may be on different scales.

SGD: choice of learning rate

Learning rate should scale as

$$\eta_t = \frac{1}{\sqrt{t}}$$

Problem: Some of the updates may be on different scales.

Solution: Let $g_{tj} = \frac{\partial L(y_t, \beta^T \mathbf{x}_t)}{\partial \beta_j}$

Scale gradients to get update rule

$$\beta_j \leftarrow \beta_j - \eta \frac{g_{tj}}{\sqrt{\sum_{s=1}^t g_{sj}^2}}$$

SGD: scaling issues

For a linear model, the SGD update is

$$\beta_j \longleftarrow \beta_j - C_t x_j$$

If x_j increases by a factor of two, the weight β_j should decrease by a factor of two.

This update doesn't respect that scaling

SGD: scaling issues

Usual solution is to “standardize” each variable — subtract out the mean and divide by the standard deviation

$$x_j \leftarrow \frac{x_j - \text{mean}(x_j)}{\sqrt{\text{var}(x_j)}}$$

But this involves “looking ahead” to compute the mean and variance, and destroys the online property of the algorithm

SGD: scaling issues

Usual solution is to “standardize” each variable — subtract out the mean and divide by the standard deviation

$$x_j \leftarrow \frac{x_j - \text{mean}(x_j)}{\sqrt{\text{var}(x_j)}}$$

But this involves “looking ahead” to compute the mean and variance, and destroys the online property of the algorithm

Solution: The mean and variance can be updated in an online manner, in constant time, by storing auxiliary variables for each component j .

SGD: Regularization

A “ridge” penalty $\lambda \sum_{j=1}^d \beta_j^2$ is easily handled.

Gradient changes by an additive term $2\lambda\beta_j$. Update becomes

$$\begin{aligned}\beta_j &\leftarrow \beta_j + \eta\{(\mathbf{y} - \pi)\mathbf{x}_j - \lambda\beta_j\} \\ &= (1 - \eta\lambda)\beta_j + \eta(\mathbf{y} - \pi)\mathbf{x}_j\end{aligned}$$

$$\beta_j \mathbf{x}_j \leftarrow (1 - \eta\lambda)\beta_j \mathbf{x}_j + \eta(\mathbf{y} - \pi)\mathbf{x}_j^2$$

Observe that this “does the right thing” whether β_j wants to be large positive or negative.

- *The penalty shrinks β_j toward zero*

What did we learn today?

- As we penalize $\|\beta\|^2$ from being too big, this *shrinks* the estimated coefficients toward zero.
- Ridge regression shrinks the coefficients. Good for high dimensions.
- If predictor variables are highly correlated, the model estimates may be unstable
- Stochastic gradient descent is a first order method that scales to large classification (and regression) problems
- Choosing the learning rate is a little tricky
- Difficult to parallelize, but not always necessary

Readings

Classification is covered in Chapter 4 of our ISL book. In particular, Section 4.3 is on logistic regression, and Section 4.4 is on linear and quadratic discriminant analysis.