

QUADRATIC LOSS: A WORKED EXAMPLE

From last class:

$$l_2(y) = (1-y)^2 = (y-1)^2 \quad m(x, y, w) = \langle x, w \rangle y \quad R_2(w) = \lambda \|w\|^2$$

$$D l_2(y) = 2(y-1) \quad D_w m(x, y, w) = x y \quad D_w R_2(w) = 2\lambda w \quad \text{Derivatives}$$

$$L(w) = \left(\sum_{(x,y)} l_2 \circ m(x, y, w) \right) + R_2(w)$$

$$D_w L(w) = \left(\sum_{(x,y)} D l_2 \circ m(x, y, w) \times D_w m(x, y, w) \right) + D_w R_2(w)$$

$$= \left(\sum_{(x,y)} 2(\langle x, w \rangle y - 1) \cdot x y \right) + 2\lambda w$$

Find the extreme value (set gradient to zero):

$$\left(\sum_{(x,y)} (\langle x, w \rangle - 1) x y \right) + \lambda w = 0$$
$$\left(\sum_{(x,y)} \langle x, w \rangle x y \right) + \lambda w = \sum x y$$

Now rearrange so we can write left side as Mw :

$$\sum y^2 x x^T w + \lambda w = \sum x y$$

$$\left(\left(\sum y^2 x x^T \right) + \lambda I \right) w = \sum x y$$

$$\left(A + \lambda I \right) w = b$$

$$M w = b$$

$$w = M^{-1} b$$

But we can do better! Determining best λ is cheap:

$$A = Q \Lambda Q^T$$

$$A + \lambda I = Q \Lambda Q^T + Q(\lambda I)Q^T$$

$$M = Q(\Lambda + \lambda I)Q^T$$

$$M^{-1} = Q(\Lambda + \lambda I)^{-1}Q^T$$

⌚ This is trivial to compute!

Behavior is intimately linked to eigenvalues and eigenvectors of A.

GRADIENT DESCENT

Let's write our loss in the format of the analysis in Goh's paper:

$$L(\omega) = \sum_{x,y} (l_2 \circ m)(x,y,m) + R_2(\omega)$$

$$= \sum_{x,y} (\langle x, \omega \rangle_y - 1)^2 + \lambda \|\omega\|^2$$

$$= \sum_{x,y} (x^T \omega y - 1)^2 + \lambda \omega^T \omega$$

$$= \sum_{x,y} (\omega^T x x^T \omega y^2 - 2 x^T \omega y + 1) + \lambda \omega^T \omega$$

$$= \omega^T \left(\sum_{x,y} x x^T y^2 \right) \omega - 2 \left(\sum_{x,y} x^T y \right) \omega + \sum_{x,y} 1 + \lambda \omega^T \omega$$

$$= \omega^T \left(\left(\sum_{x,y} x x^T y^2 \right) + \lambda I \right) \omega - \left(2 \sum_{x,y} x^T y \right) \omega + \sum 1$$

$$f(\omega) = \frac{1}{2} \omega^T A \omega - b^T \omega$$

$$\nabla_{\omega} f(\omega) = A \omega - b$$

$$\nabla_{\omega} f = 0$$

$$A \omega^* = b$$

$$\omega^* = A^{-1} b$$

(we will use this soon)

GRADIENT ITERATION

$$\omega^{k+1} = \omega^k - \alpha(A\omega^k - b)$$

A is symmetric, so admits eigenanalysis:

$$A = Q\Lambda Q^T$$

Now let's look at consecutive iterations:

$$\omega^{k+1} = \omega^k - \alpha(A\omega^k - b)$$

$$\omega^{k+1} = (I - \alpha A)\omega^k + \alpha b$$

$$\begin{aligned} I - \alpha A &= QQ^T - \alpha Q\Lambda Q^T \\ &= Q(I - \alpha\Lambda)Q^T \end{aligned}$$

$$\begin{aligned} (I - \alpha A)^2 &= Q(I - \alpha\Lambda)Q^T Q(I - \alpha\Lambda)Q^T \\ &= Q(I - \alpha\Lambda)^2 Q^T \end{aligned}$$

$I - \alpha\Lambda$ is a simpler expression than $I - \alpha A$: the former is a diagonal matrix. Let's try to get gradient descent to look like $I - \alpha\Lambda$. First, let's translate the whole problem so that the minimum is at zero, and then let's write everything in terms of vectors multiplied by Q^T :

$$\omega^* = A^{-1}b$$

$$\omega^* = Q\Lambda^{-1}Q^T b$$

$$Q^T \omega^* = \Lambda^{-1} Q^T b$$

$$\Lambda Q^T \omega^* = Q^T b$$

$$x^k = Q^T (\omega^k - \omega^*)$$

$$Qx^k = \omega^k - \omega^*$$

$$\omega^k = Qx^k + \omega^*$$

$$\omega^{k+1} = \omega^k - \alpha (A\omega^k - b)$$

$$Q^T \omega^{k+1} = Q^T \omega^k - \alpha Q^T (A\omega^k - b)$$

$$Q^T \omega^{k+1} - Q^T \omega^* = Q^T \omega^k - Q^T \omega^* - \alpha Q^T (A\omega^k - b)$$

$$x^{k+1} = x^k - \alpha (Q^T Q \Lambda Q^T \omega^k - Q^T b)$$

$$x^{k+1} = x^k - \alpha (\Lambda Q^T \omega^k - \Lambda Q^T \omega^* + \Lambda Q^T \omega^* - Q^T b)$$

$$x^{k+1} = x^k - \alpha \Lambda x^k$$

$$x^{k+1} = (I - \alpha \Lambda) x^k$$

$$x^k = (I - \alpha \Lambda)^k x^0$$

$$Qx^k = Q(I - \alpha \Lambda)^k x^0$$

$$Qx^k + \omega^* = Q(I - \alpha \Lambda)^k x^0 + \omega^*$$

$$\omega^k = Q(I - \alpha \Lambda)^k x^0 + \omega^*$$

$$\omega^k = Q(I - \alpha \Lambda)^k Q^T (\omega^k - \omega^*) + \omega^*$$

How do we make sure that ω^k will go to ω^* ?

The error $\omega^k - \omega^*$ must go to zero. In addition, each coordinate of $Q^T(\omega^k - \omega^*)$ goes to zero (or fails to) independently of the other: each gets scaled by $(I - \alpha \Lambda)$ every iteration.

This is the true picture of gradient descent: the error vector is multidimensional, and the amount it shrinks on coordinate i (of the eigenspace) is given by $|1 - \alpha \lambda_i|$, and each coordinate shrinks independently of the other.

Given an initial guess x^0 , the magnitude of the (quadratic) error is:

$$\sum (x_i^0)^2$$

Even if each x_i is identical, each iteration reduces each coordinate proportionally to λ_i , so the features corresponding to smaller eigenvalues will take longer to converge to their minimum. The difference between early iterations and late iterations, then, is due to features coming from small eigenvalues.

WHAT ABOUT REGULARIZATION?

Intuitively, regularization should make optimization want smaller w vectors,

$$L(w) = \sum_{x,y} \text{loss-of-sample}(y, \text{margin}(x,y,w)) + \frac{\lambda \|w\|^2}{2}$$

This becomes directly obvious when we consider a gradient descent update:

$$\nabla_w L(w) = (\nabla_w \sum \text{loss-of-sample}) + \lambda w$$

So a step taken with a regularized loss is like first moving a little towards zero, and then taking a regular step:

$$\begin{aligned} w^{k+1} &= w^k - \alpha \nabla_w L(w) \\ &= w^k - \alpha ((\nabla_w \sum \text{loss-of-sample}) + \lambda w^k) \\ &= w^k - \alpha (\nabla_w \sum \text{loss-of-sample}) - \alpha \lambda w^k \\ &= (1 - \alpha \lambda) w^k - \alpha (\nabla_w \sum \text{loss-of-sample}) \end{aligned}$$

(The notation here gets a little confusing: λ is the regularization term, and λ_i are the eigenvalues of A)

Finally, let's consider the difference between the solution of the regularized problem and the solution of the unregularized problem. The difference, incredibly, is analogous to that of choosing whether to stop gradient descent early!

$$\begin{aligned} 2 f(w) &= w^T A w - 2w^T b + \lambda \|w\|^2 \\ 2 \nabla f &= 2A w - 2b + 2\lambda w \\ \nabla f &= (A + \lambda I) w - b \end{aligned}$$

$$(A + \lambda I) w = b$$

$$\begin{aligned} w^* &= A^{-1} b \\ w'^* &= (A + \lambda I)^{-1} b \\ w^* - w'^* &= (A^{-1} - (A + \lambda I)^{-1}) b \end{aligned}$$

But recall our eigenanalysis:

$$A = Q \Lambda Q^T$$

$$I = Q Q^T$$

$$\begin{aligned} w^* - w'^* &= ((Q \Lambda Q^T)^{-1} - (Q \Lambda Q^T + \lambda Q Q^T)^{-1}) b \\ &= ((Q \Lambda Q^T)^{-1} - (Q (\Lambda + \lambda I) Q^T)^{-1}) b \end{aligned}$$

Stupid linear algebra tricks, now:

1) If Q is orthogonal, $Q^{-1} = Q^T$

2) $(AB)^T = B^T A^T$

$$\begin{aligned}(Q \Lambda Q^T)^{-1} &= (Q^T)^{-1} (\Lambda)^{-1} \\ &= Q \Lambda^{-1} Q^T \\ &= Q \Lambda^{-1} Q^T \quad (!)\end{aligned}$$

$$\begin{aligned}\omega^* - \omega'^* &= \left(Q \operatorname{diag} \left(\frac{1}{\lambda_i} \right) Q^T - Q \operatorname{diag} \left(\frac{1}{\lambda + \lambda_i} \right) Q^T \right) b \\ &= Q \operatorname{diag} \left(\frac{\lambda}{\lambda + \lambda_i}, \frac{1}{\lambda_i} \right) Q^T b \quad \left(\frac{1}{\lambda} - \frac{1}{\lambda + \lambda_i} \right)\end{aligned}$$

What happens when $\lambda \gg \lambda_i$?

What happens when $\lambda_i \gg \lambda$?