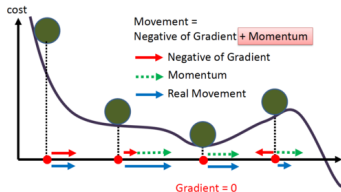


Optimization in Machine Learning

First order methods: GD on quadratic forms

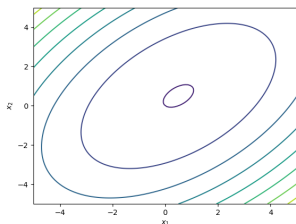


Learning goals

- Definition
- Max. Likelihood
- Normal regression
- Risk Minimization

GD ON QUADRATIC FORMS

- We consider the quadratic function $q(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$.
- We assume that the Hessian $\mathbf{H} = \mathbf{A}$ is symmetric and invertible
- The optimal solution is $\mathbf{x}^* = \mathbf{A}^{-1} \mathbf{b}$
- As $\nabla q(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}$, the iterations of gradient descent are
$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha(\mathbf{A} \mathbf{x}^{[t]} - \mathbf{b})$$



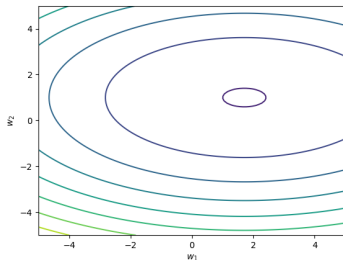
The following slides follow the blogpost by Goh, "Why Momentum Really Works", Distill, 2017. <http://doi.org/10.23915/distill.00006>

GD ON QUADRATIC FORMS

For \mathbf{A} , there exists an eigenvalue decomposition:

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$$

where the columns of \mathbf{V} contain the eigenvectors \mathbf{v}_i and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ contains the eigenvalues λ_i sorted from smallest to biggest. We perform a change of basis $\mathbf{w}^{[t]} = \mathbf{V}^\top (\mathbf{x}^{[t]} - \mathbf{x}^*)$ to its eigenspace, where all dimensions act independently.



GD ON QUADRATIC FORMS

We get

$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha \left(\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b} \right)$$

$$\mathbf{w}^{[t]} = \mathbf{V}^\top (\mathbf{x}^{[t]} - \mathbf{x}^*)$$

$$\mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{x}^* = \mathbf{x}^{[t]}$$

and with $\mathbf{x}^{[t]} = \mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{x}^*$ we can write a GD step as:

GD ON QUADRATIC FORMS

$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha (\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

$$\begin{array}{l|l} \mathbf{V} \cdot \mathbf{w}^{[t+1]} + \mathbf{x}^* = \mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{x}^* - \alpha (\mathbf{A} \cdot \mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{A} \cdot \mathbf{x}^* - \mathbf{b}) & - \mathbf{x}^* \\ \mathbf{V} \cdot \mathbf{w}^{[t+1]} = \mathbf{V} \cdot \mathbf{w}^{[t]} - \alpha (\mathbf{A} \cdot \mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{A}\mathbf{x}^* - \mathbf{b}) & \mathbf{V}^\top \text{ (NB: } \mathbf{V}^\top \mathbf{V} = \mathbf{I} \text{)} \\ \mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - \alpha \mathbf{V}^\top (\mathbf{A}\mathbf{V} \cdot \mathbf{w}^{[t]} + \mathbf{A}\mathbf{x}^* - \mathbf{b}) & \mathbf{A}\mathbf{x}^* - \mathbf{b} = \mathbf{0} \\ \mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - \alpha \mathbf{V}^\top \mathbf{A}\mathbf{V} \cdot \mathbf{w}^{[t]} & \mathbf{V}^\top \mathbf{A}\mathbf{V} = \Sigma \\ \mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - \alpha \Sigma \mathbf{w}^{[t]} & \end{array}$$

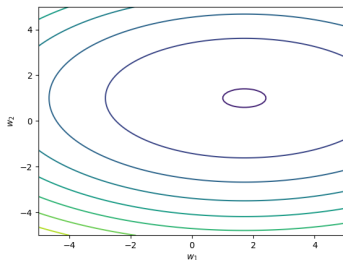
Which means:

$$w_i^{[t+1]} = w_i^{[t]} - \alpha \lambda_i w_i^{[t]} = (1 - \alpha \lambda_i) w_i^{[t]} = (1 - \alpha \lambda_i)^{t+1} w_i^{[0]}$$

GD ON QUADRATIC FORMS

If we now perform GD on \mathbf{w} , we get

$$\begin{aligned}w_i^{[t+1]} &= w_i^{[t]} - \alpha \lambda_i w_i^{[t]} \\&= (1 - \alpha \lambda_i) w_i^{[t]} = (1 - \alpha \lambda_i)^{t+1} w_i^{[0]}\end{aligned}$$



GD ON QUADRATIC FORMS

Moving back to the original space, we get

$$\mathbf{x}^{[t]} - \mathbf{x}^* = \mathbf{V} \cdot \mathbf{w}^{[t]} = \sum_{i=1}^d w_i^{[0]} (1 - \alpha \lambda_i)^t \mathbf{v}_i$$

This allows a very intuitive interpretation: each element of $w^{[0]}$ is the component of the error in the initial guess in the eigenbasis and decays with a rate of $1 - \alpha \lambda_i$.

For most step sizes, the eigenvectors with the largest eigenvalues converge the fastest.

GD ON QUADRATIC FORMS

We now consider the contribution of each eigenvector to the total loss

$$q(\mathbf{x}^{[t]}) - q(\mathbf{x}^*) = \frac{1}{2} \sum_i^d (1 - \alpha \lambda_i)^{2t} \lambda_i (w_i^{[0]})^2$$

