

ECE/CS 559 Lecture 8 Th 9/9

Last time: Linear Regression with Squared Loss

- Data: $\{(x,y), \dots\}$ $x \in \mathbb{R}^n$ $y \in \mathbb{R}$
- Predictor: $f(x; w) = w^T x$ $w \in \mathbb{R}^n$
- Goal: minimize $R(w) = \sum_{(x,y) \in \text{Data}} \|y - w^T x\|^2$
 $\ell(y, w^T x) \equiv \ell((x,y), w)$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \cdot \begin{bmatrix} 1 & x_1 & \dots & x_n \end{bmatrix} = \begin{bmatrix} 1 & y \end{bmatrix}$$

it's a quadratic equation

$$w = Y X^T (X X^T)^{-1}$$

Pseudo inverse



1-D Example: $R(w) = w^4$ $\frac{dR}{dw} = 4w^3$ $\eta = \frac{1}{8}$

$$w(0) = 1 \quad w(1) \leftarrow w(0) - \nabla R(w(0)) \cdot \eta = 1 - 4 \cdot 1^3 \cdot \frac{1}{8} = \frac{1}{2}$$

$$w(2) \leftarrow w(1) - \nabla R(w(1)) \cdot \eta = \frac{1}{2} - 4 \left(\frac{1}{2}\right)^3 \cdot \frac{1}{8} = \frac{7}{16}$$

• What's happening?

$$w(t) \leftarrow w(t-1) - 4 w(t-1)^3 \eta = [1 - 4\eta w(t-1)^2] w(t-1)$$

- If η is small enough, always diminishes. (< 1) since bounded from below (by 0) will converge.
- But will it converge to 0? Yes, think about why.
- What if η is too big? Try $\eta = \frac{1}{2}$ with $w(0) = 1$.

③ Delta Rule

• Squared loss for single neuron with differentiable $\phi(\cdot)$:

$$R(w) = \sum_{(x,y)} (y - \phi(w^T x))^2$$

$$\nabla R(w) = -2 \sum_{(x,y)} (y - \phi(w^T x)) \phi'(w^T x) x^T$$

chain rule

$$w \leftarrow w + 2\eta \sum_{(x,y)} (y - \phi(w^T x)) \phi'(w^T x) x$$

• 1 data point at a time:

$$w \leftarrow w + \eta (y - \phi(w^T x)) \phi'(w^T x) x$$

① Gradient Descent

Goal: minimize $R(w) = \sum_{(x,y) \in \text{Data}} \ell(y, f(x; w)) = \sum_{(x,y) \in \text{Data}} \|y - f(x; w)\|^2$

Idea: Set the derivative to zero!

Gradient: $\nabla R = \left[\frac{\partial R}{\partial w_1}, \dots, \frac{\partial R}{\partial w_n} \right]$ reduce step-by-step.

Each step: Change w by Δw : $w' \leftarrow w + \Delta w$

$$R(w') = R(w + \Delta w) = R(w) + \nabla R(w)^T \Delta w + O(\|\Delta w\|^2)$$

How should we choose Δw to get the best reduction? Taylor expansion ignore

- To minimize $\nabla R^T \Delta w$, let $\Delta w \propto -\nabla R$: $w' \leftarrow w - \eta \nabla R(w)$
- proportional

② Widrow-Hoff LMS Algorithm

• Let's apply gradient descent to linear regression. Recall:

$$\nabla R(w) = -2 \sum_{(x,y) \in \text{Data}} (y - w^T x) x^T$$

same shape as w

$$w \leftarrow w - \eta \nabla R(w) = w + 2\eta \sum_{(x,y)} (y - w^T x) x^T$$

absorb

• When $y \in \mathbb{R}^1$ then $w = w^T$ and $y - w^T x \in \mathbb{R}$ The updates become: $w \leftarrow w + \eta \sum_{(x,y)} (y - w^T x) x$

• We could do this 1 data point at a time:

For each $(x,y) \in \text{Data}$: $w \leftarrow w + \eta (y - w^T x) x$

Similar to the perceptron learning algorithm!

Example $\phi(\cdot)$:

• Sigmoid: $\phi(v) = \frac{1}{1 + e^{-av}}$ Since $1 - \phi(v) = \frac{e^{-av}}{1 + e^{-av}}$

$$\phi'(v) = \frac{0 - (-ae^{-av})}{(1 + e^{-av})^2} = \frac{ae^{-av}}{(1 + e^{-av})^2} = a\phi(v)(1 - \phi(v))$$

• ReLU: $\phi(v) = \begin{cases} v & v \geq 0 \\ 0 & \text{otherwise} \end{cases}$

$$\phi'(v) = \begin{cases} 1 & v \geq 0 \\ 0 & \text{otherwise} \end{cases} = \text{step}(v)$$