

ECE/CS 559 Lecture 7

9/17

Last time : Perception Learning Example Supervised Learning

- Data: $\{(x_i, y_i), \dots\}$ $x \in \mathcal{X}$ (feature space), $y \in \mathcal{Y}$ (label space)
- Task: Predict/initiate y using only x .
- How: Use a predictor neural network $y_{\text{pred}}(x) = f(x; w)$
- Loss: how far y and $f(x; w)$ are: $\ell(y, f) = \frac{1}{2} \{y - f\}^2$
- Risk: (empirical) $R(w) = \frac{1}{|\text{Data}|} \sum_{(x_i, y_i) \in \text{Data}} \ell(y_i, f(x_i; w))$
- Goal: Minimize (empirical) risk,
make few mistakes/errors / stay close to y .

• Types of algorithms:

- See (x_i, y_i) , update w
- See (x_i, y_i) , update w
- See (x_i, y_i) , update w

- Get Data = $\{(x_i, y_i), \dots\}$
- Epoch 1: use Data, update w
- Epoch 2: use Data, update w



① Linear Regression with Squared Loss

- Data: $\{(x_i, y_i), \dots\} \quad x \in \mathbb{R}^n \quad y \in \mathbb{R}^m$
- Predictor: $f(x; w) = Wx \quad W \in \mathbb{R}^{m \times n}$
(this is a neuron with $g(v) = v$, identity activation)
- Goal: minimize $R(w) = \sum_{(x_i, y_i) \in \text{Data}} \|y_i - Wx_i\|^2$

Analytic Solution: Optimality conditions, unconstrained $\Rightarrow \nabla_w R = 0$

$$\begin{aligned} \text{• Single data point: } l &= \sum_{k=1}^m (y_k - \sum_{j=1}^n w_{kj} x_j)^2 \quad \text{gradient} \\ \frac{\partial l}{\partial w_{kj}} &= -2x_j (y_k - \sum_{j=1}^n w_{kj} x_j) = -2y_k x_j + \sum_{j=1}^n w_{kj} (x_j x_j) \\ \downarrow \begin{matrix} i \\ m \end{matrix} \quad \rightarrow \quad \nabla_w l &= -2y x^T + W x x^T \end{aligned}$$

$$\text{• All data points: } R(w) = \sum_{(x_i, y_i) \in \text{Data}} \ell(y_i, Wx_i) \Rightarrow \nabla_w R = \sum_{(x_i, y_i) \in \text{Data}} \nabla_l$$

$$\begin{aligned} \nabla_w R &= \sum_{(x_i, y_i) \in \text{Data}} -2y x^T + 2W \sum_{x \in \text{Data}} x x^T \\ &= -2 \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}}_{|\text{Data}|} \underbrace{\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^T}_{\text{Data}} + 2W \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}}_{\text{Data}} \underbrace{\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}^T}_{\text{Data}} \\ &= -2Y X^T + 2W X X^T \end{aligned}$$

$$\bullet \nabla_w R = 0 \Rightarrow -2Y X^T + 2W X X^T = 0$$

$$\Rightarrow W = Y X^T (X X^T)^{-1} \quad \text{Moore-Penrose Pseudo inverse}$$

$$\text{Compare to: } W X \approx Y \quad \text{not invertible if } |\text{Data}| > n \\ \text{• } [W] \begin{bmatrix} \vdots & \vdots & \dots & \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \vdots & \vdots & \dots & \vdots \end{bmatrix} \quad \text{but it's "like" we inverted if above.}$$

• Where are the updates? One single batch update!

Intuition: it's a quadratic equation



• Other shapes → solve by iteration/updates

② Gradient Descent

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

Idea: Set the derivative to zero!

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

$$\text{Each step: Change } w \text{ by } \Delta w: \quad w' \leftarrow w + \Delta w.$$

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

Taylor expansion ignore

How should we choose Δw to get the best reduction?

$$\text{To minimize } \nabla_w R(w)^T \Delta w, \text{ let } \Delta w \propto -\nabla_w R(w): \quad w' \leftarrow w - \eta \nabla_w R(w)$$

proportional

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

$$\begin{aligned} w(0) &= 1 & w(1) &\leftarrow w(0) - \nabla R(w(0)) \cdot \eta \\ && 1 &- 4 \cdot \frac{1}{8} = \frac{1}{2} \\ w(2) &\leftarrow w(1) - \nabla R(w(1)) \cdot \eta \\ &\frac{1}{2} &- 4 \left(\frac{1}{2} \right)^3 \cdot \frac{1}{8} = \frac{3}{16} \end{aligned}$$

• 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 $w(0) = 1$.