Introduction to Machine Learning

Supervised Regression Deep Dive: Proof OLS Regression

× × 0 × × ×



Learning goals

 Understand analytical derivation of OLS estimator for LM

ANALYTICAL OPTIMIZATION

• Special property of LM with L2 loss: analytical solution available

$$\begin{split} \hat{\boldsymbol{\theta}} \in \mathop{\arg\min}_{\boldsymbol{\theta}} \mathcal{R}_{emp}(\boldsymbol{\theta}) &= \mathop{\arg\min}_{\boldsymbol{\theta}} \sum_{i=1}^{n} \left(\boldsymbol{y}^{(i)} - \boldsymbol{\theta}^{\top} \mathbf{x}^{(i)} \right)^{2} \\ &= \mathop{\arg\min}_{\boldsymbol{\theta}} \| \mathbf{y} - \mathbf{X} \boldsymbol{\theta} \|_{2}^{2} \end{split}$$

• Find via normal equations

$$\frac{\partial \mathcal{R}_{\mathsf{emp}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{0}$$

• Solution: ordinary-least-squares (OLS) estimator

$$\hat{oldsymbol{ heta}} = (\mathbf{X}^{ op} \mathbf{X})^{-1} \mathbf{X}^{ op} \mathbf{y}$$

ANALYTICAL OPTIMIZATION – PROOF

$$\mathcal{R}_{emp}(\boldsymbol{\theta}) = \sum_{i=1}^{n} \left(\underbrace{\mathbf{y}^{(i)} - \boldsymbol{\theta}^{\top} \mathbf{x}^{(i)}}_{=:\epsilon_{i}} \right)^{2} = \| \underbrace{\mathbf{y} - \mathbf{X}\boldsymbol{\theta}}_{=:\epsilon} \|_{2}^{2}; \quad \boldsymbol{\theta} \in \mathbb{R}^{\tilde{p}} \text{ with } \tilde{p} := p + 1$$

$$\begin{array}{rcl} 0 & = & \frac{\partial \mathcal{R}_{emp}(\theta)}{\partial \theta} \mbox{ (sum notation)} & 0 & = & \frac{\partial \mathcal{R}_{emp}(\theta)}{\partial \theta} \mbox{ (matrix notation)} \\ 0 & = & \frac{\partial}{\partial \theta} \sum_{i=1}^{n} \epsilon_i^2 \mbox{ | sum \& chain rule} & 0 & = & \frac{\partial \|\epsilon\|_2^2}{\partial \theta} \\ 0 & = & \sum_{i=1}^{n} \frac{\partial \epsilon_i^2}{\partial \epsilon_i} \frac{\partial \epsilon_i}{\partial \theta} & 0 & = & \frac{\partial \epsilon^{\top} \epsilon}{\partial \theta} \mbox{ | chain rule} \\ 0 & = & \sum_{i=1}^{n} \frac{\partial \epsilon_i^2}{\partial \theta} \frac{\partial \epsilon_i}{\partial \theta} & 0 & = & \frac{\partial \epsilon^{\top} \epsilon}{\partial \theta} \mbox{ | chain rule} \\ 0 & = & \sum_{i=1}^{n} 2\epsilon_i(-1)(\mathbf{x}^{(i)})^{\top} & 0 & = & \frac{\partial \epsilon^{\top} \epsilon}{\partial \epsilon} \cdot \frac{\partial \epsilon}{\partial \theta} \\ 0 & = & \sum_{i=1}^{n} (\mathbf{y}^{(i)} - \theta^{\top} \mathbf{x}^{(i)})(\mathbf{x}^{(i)})^{\top} & 0 & = & (\mathbf{y} - \mathbf{X}\theta)^{\top} \mathbf{x} \\ 0 & = & \sum_{i=1}^{n} (\mathbf{y}^{(i)} - \theta^{\top} \mathbf{x}^{(i)})(\mathbf{x}^{(i)})^{\top} & 0 & = & \mathbf{y}^{\top} \mathbf{x} - \theta^{\top} \mathbf{x}^{\top} \mathbf{x} \\ \theta^{\top} \sum_{i=1}^{n} \mathbf{x}^{(i)}(\mathbf{x}^{(i)})^{\top} & = & \sum_{i=1}^{n} \mathbf{y}^{(i)}(\mathbf{x}^{(i)})^{\top} \mbox{ | transpose} & \theta^{\top} \mathbf{x}^{\top} \mathbf{x} & = & \mathbf{y}^{\top} \mathbf{x} \mbox{ | transpose} \\ \left(\sum_{i=1}^{n} \frac{\mathbf{x}^{(i)}(\mathbf{x}^{(i)})^{\top}}{\hat{p} \times \hat{p}} \mbox{ | ternspose} & \mathbf{x}^{\top} \mathbf{x} \theta & = & \mathbf{x}^{\top} \mathbf{y} \\ \left(\sum_{i=1}^{n} \mathbf{x}^{(i)}(\mathbf{x}^{(i)})^{\top} \mbox{ | ternspose} & \mathbf{x}^{\top} \mathbf{x} \mbox{ | ternspose} \\ \mathbf{x}^{\top} \mathbf{x} \theta & = & \mathbf{x}^{\top} \mathbf{y} \\ (\mathbf{x}^{\top} \mathbf{x}) \theta & = & \mathbf{x}^{\top} \mathbf{y} \\ NB: & \sum_{i=1}^{n} \mathbf{x}^{(i)}(\mathbf{x}^{(i)})^{\top} \mbox{ | seasy to show (try it!) - and good to remember (this is basically the estimation of Cov(X)) \\ \end{array}$$

× × 0 × × ×