

Regularization

Instructor: Jessica Wu -- Harvey Mudd College

The instructor gratefully acknowledges Andrew Ng (Stanford), Eric Eaton (UPenn), David Kauchak (Pomona), and the many others who made their course materials freely available online.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Regularization

Learning Goals

- Describe goal of regularization
- Describe how to solve regularized regression using gradient descent and normal equations
- Describe common regularizers and tradeoffs









$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^{n} \left(h_{\boldsymbol{\theta}} \left(\boldsymbol{x}^{(i)} \right) - \boldsymbol{y}^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} \theta_j^2$$

• Note that $\sum_{j=1}^{d} \theta_{j}^{2} = \| \theta_{1:d} \|_{2}^{2}$

- this is magnitude of the feature coefficient vector

- We can also think of this as $\sum_{j=1}^{d} (\theta_j 0)^2 = \|\boldsymbol{\theta}_{1:d} \mathbf{0}\|_2^2$
 - L_2 -regularization pulls coefficients towards 0



Regularized Linear Regression

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^{n} \left(h_{\boldsymbol{\theta}} \left(\boldsymbol{x}^{(i)} \right) - \boldsymbol{y}^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} \theta_j^2$$

Gradient Update (single training example)

$$\frac{\partial}{\partial \theta_0} J(\boldsymbol{\theta})$$
$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$



Regularized Linear Regression $J(\boldsymbol{\theta}) = \frac{1}{2} \left[(\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})^T (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y}) + \lambda \boldsymbol{\theta}^T \boldsymbol{\theta} \right]$ really only $\boldsymbol{\theta}_{1:d}$ Closed-Form Solution $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = 0$

Based on slide by Eric Eaton



Generalizing Regularization

Regularized Linear Regression

$$I(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^{n} \left(h_{\boldsymbol{\theta}} \left(\boldsymbol{x}^{(i)} \right) - \boldsymbol{y}^{(i)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{d} \theta_j^2$$

Generally

$$R_n(\boldsymbol{\theta}) = \|\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y}\|_2^2$$
$$J_{n,\lambda}(\boldsymbol{\theta}) = R_n(\boldsymbol{\theta}) + \lambda z(\boldsymbol{\theta}) \qquad z(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_2^2$$
empirical risk regularization model fit to data

- $z(\boldsymbol{\theta})$ regularization
- biases parameters toward default values (e.g. 0)
- resists setting parameters away from zero when data (weakly) tells us otherwise



Common Regularizers

Claim: If f and g are convex functions, then so if f + g. Claim: p-norms are convex for $p \ge 1$.

Common names for regularized regression

| name | loss | regularization |
|------------------------------|----------|-----------------|
| ordinary least squares (OLS) | squared | none |
| ridge regression | squared | L ₂ |
| lasso regression | squared | L ₁ |
| elastic regression | squared | $L_{1} + L_{2}$ |
| logistic regression | logistic | |



