Lecture 7 - Overfitting and Regularization

Dan Sheldon

September 26, 2012

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

- What is Overfitting?
- How to Diagnose Overfitting

(ロ)、(型)、(E)、(E)、 E) の(の)

Regularization

What is Overfitting?

Demo: polynomials



What is Overfitting?

Complex decision boundaries



ł

Overfitting is learning a model that fits the training data very well, but does not *generalize* well.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

(Generalize well = predict accurately for new examples.)

How to Diagnose Overfitting?

Exercise

How to Diagnose Overfitting?

Exercise

Reserve some data to test whether hypothesis generalizes well



Train Data vs. Test Data

Very important (and simple) methodology

Start with N training examples

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_N, y_N)$$

- Split randomly into train and test sets
- To fit the model, minimize cost on train data only

$$J_{\mathsf{train}}(\mathbf{w}) = \sum_{i \in \mathsf{train}} \mathsf{cost}(h(\mathbf{x}_i), y_i)$$

To evaluate the fit, measure cost on test set

$$J_{\mathsf{test}}(\mathbf{w}) = \sum_{i \in \mathsf{test}} \mathsf{cost}(h(\mathbf{x}_i), y_i)$$

How to Diagnose Overfitting?

Example: cost function vs. degree of polynomial



▲ 臣 ▶ 臣 • • • • • •

How to Diagnose Overfitting?

$$(X_{1}, X_{2}) \longmapsto (X_{1}, X_{2}, X_{1}^{*}, X_{2}^{*}, X_{1}X_{2})$$

Example: cost function vs. number of features in book data



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

owers Cost vs. Complexity .10 1,2,3,... General phenomenon: training/test cost vs. model "complexity" traihing $\lambda = 100$ (omplex, 'ty)

What Makes a Model Complex?

- Polynomial: higher degree
- Book data: more features

• Linear functions $(h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x})$: large weights

$$h(x) = w_{0} + w_{1}x_{1} + w_{2}x_{2} + \dots + w_{d}x_{d}$$

$$h(x) = w_{0} + w_{1}x + w_{2}x^{2} + \dots + w_{d}x^{d}$$

$$\int_{V_{4}=0}^{T} \frac{1}{w_{d}=0}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Large Weights



Which is more complex?

$$y = -3.94 + 0.18x_1 + .34x_2$$

VS.

$$y = 2842 - 957x_1 + 300x_2$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



Intuition: large weights \rightarrow high complexity

So, modify the cost function to penalize large weights. For linear regression, the new cost function is:

$$J(\mathbf{w}) = \lambda \sum_{j=0}^{d} w_j^2 + \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$
New term $\rightarrow penalty on high weights$

 λ controls trade-off between model complexity and fit

$$\chi = .01 \rightarrow$$
 more emphasis on tit
 $\chi = 1000 \rightarrow$ more empasis on weights

Discussion

Regularization is really important!!! Why?



Derivation (review on your own)

$$J(\mathbf{w}) = \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \sum_{j=0}^{d} w_j^2$$
$$= (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}.$$

Set derivative to zero

$$0 = \frac{d}{d\mathbf{w}}J(\mathbf{w})$$

$$0 = 2(X\mathbf{w} - \mathbf{y})^T X + 2\lambda \mathbf{w}^T$$

$$0 = X^T (X\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}$$

$$X^T X \mathbf{w} + \lambda \mathbf{w} = X^T \mathbf{y}$$

$$X^T X + \lambda I) \mathbf{w} = X^T \mathbf{y}$$

$$\mathbf{w} = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへぐ

Regularized Gradient Descent for Linear Regression

$$J(\mathbf{w}) = \lambda \sum_{j=0}^{d} w_j^2 + \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

Repeat until convergence

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w}), \qquad j = 0, \dots, d.$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Regularized Gradient Descent for Linear Regression

$$J(\mathbf{w}) = \lambda \sum_{j=0}^{d} w_j^2 + \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Repeat until convergence

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w}), \qquad j = 0, \dots, d.$$

$$w_j \leftarrow w_j - \alpha \left(2\lambda w_j + 2\sum_{i=1}^N (h_{\mathbf{w}}(\mathbf{x}_i) - y_i) x_{i,j} \right)$$

Regularized Gradient Descent for Linear Regression

$$J(\mathbf{w}) = \lambda \sum_{j=0}^{d} w_j^2 + \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$$

Repeat until convergence

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w}), \qquad j = 0, \dots, d.$$

$$w_{j} \leftarrow w_{j} - \alpha \left(2\lambda w_{j} + 2\sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_{i}) - y_{i})x_{i,j} \right)$$

$$w_{j} \leftarrow w_{j}(1 - 2\lambda\alpha) - 2\alpha \sum_{i=1}^{N} (h_{\mathbf{w}}(\mathbf{x}_{i}) - y_{i})x_{i,j}$$

$$sh_{f} h_{k} h_{k} g \qquad \text{old } f \text{ radies for } step$$

Regularized Gradient Descent for Logistic Regression

$$J(\mathbf{w}) = \lambda \sum_{j=0}^{d} w_j^2 + \sum_{i=1}^{N \left(-\gamma_i \log h_i(x_i) - (1-\gamma_i)\log(1-h(x_i))\right)} \frac{1}{(1-h_i(x_i))}$$

×

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Repeat until convergence. For $j = 0, \ldots, d$:

What You Need To Know

- Concept of overfitting
- Diagnosis: train/test sets
- Regularized cost function (penalize weights)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Regularized gradient descent
- See it work