

Linear Regression: Vectorization, Regularization

These slides were assembled by Byron Boots, with grateful acknowledgement to Eric Eaton and the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Last Time: Basis Functions

 $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{i=1}^{n} \theta_{j} x_{j}$

j=0

• Basic linear model:

• More general linear model: $h_{\theta}(\boldsymbol{x}) = \sum_{j=0}^{\infty} \theta_j \phi_j(\boldsymbol{x})$

- Once we have replaced the data by the outputs of the basis functions, fitting the generalized model is exactly the same problem as fitting the basic model
 - Unless we use the kernel trick more on that when we cover support vector machines

Vectorization

- Benefits of vectorization
 - More compact equations
 - Faster code (using optimized matrix libraries)
- Consider our model:

$$h(\boldsymbol{x}) = \sum_{j=0}^{n} \theta_j x_j$$

d

• Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \boldsymbol{x}^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

• Can write the model in vectorized form as $h(\boldsymbol{x}) = \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$

Vectorization

Consider our model for n instances:



• Can write the model in vectorized form as $\,h_{oldsymbol{ heta}}(oldsymbol{x}) = oldsymbol{X}oldsymbol{ heta}$

Vectorization

• For the linear regression cost function:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left(h_{\boldsymbol{\theta}} \left(\boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^{2}$$
$$= \frac{1}{2n} \sum_{i=1}^{n} \left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}^{(i)} - y^{(i)} \right)^{2}$$
$$\overset{\mathbb{R}^{n \times (d+1)}}{\overset{\mathbb{R}^{n \times (d+1)}}$$

 $\mathbb{R}^{\tilde{n} \times 1}$

 $=\frac{1}{2n}\left(\boldsymbol{X}\boldsymbol{\theta}-\boldsymbol{y}\right)^{\mathsf{T}}\left(\boldsymbol{X}\boldsymbol{\theta}-\boldsymbol{y}\right)$

 \mathbb{R}

 $1 \times n$

Let:

$$\boldsymbol{y} = \left[egin{array}{c} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{array}
ight]$$

Closed Form Solution

• Instead of using GD, solve for optimal heta analytically

– Notice that the solution is when $\frac{\partial}{\partial \theta} J(\theta) = 0$

• Derivation:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2n} \left(\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y} \right)^{\mathsf{T}} \left(\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y} \right)^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} - \boldsymbol{y}^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} - \boldsymbol{y}^{\mathsf{T}} \mathbf{X}\boldsymbol{\theta} \\ \propto \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y} \\ \propto \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - 2\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y}$$

Take derivative and set equal to 0, then solve for $\boldsymbol{\theta}$: $\frac{\partial}{\partial \boldsymbol{\theta}} \left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} \boldsymbol{\theta} - 2 \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y} \right) = 0$ $(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}) \boldsymbol{\theta} - \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} = 0$ $(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}) \boldsymbol{\theta} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y}$

Closed Form Solution:

$$\boldsymbol{\theta} = (\boldsymbol{X}^{\intercal}\boldsymbol{X})^{-1}\boldsymbol{X}^{\intercal}\boldsymbol{y}$$

Closed Form Solution

• Can obtain heta by simply plugging X and y into

- If $X^T X$ is not invertible (i.e., singular), may need to:
 - Use pseudo-inverse instead of the inverse
 - In python, numpy.linalg.pinv(a)
 - Remove redundant (not linearly independent) features
 - Remove extra features to ensure that $d \le n$

Gradient Descent vs Closed Form

Gradient Descent Closed Form Solution Requires multiple iterations Non-iterative • Need to choose α No need for α • Works well when n is large Slow if n is large \bullet • - Computing $(X^{T}X)^{-1}$ is Can support incremental ۲ roughly $O(n^3)$ learning

Improving Learning: Feature Scaling

• Idea: Ensure that feature have similar scales



• Makes gradient descent converge much faster

Feature Standardization

- Rescales features to have zero mean and unit variance
 - Let μ_j be the mean of feature j:

$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$$

– Replace each value with:

$$x_j^{(i)} \leftarrow rac{x_j^{(i)} - \mu_j}{s_j}$$
 for j = 1...c (not x₀!)

- s_j is the standard deviation of feature j
- Could also use the range of feature j (max_i min_i) for s_i
- Must apply the same transformation to instances for both training and prediction
- Outliers can cause problems



Overfitting:

- The learned hypothesis may fit the training set very well ($J(\pmb{\theta}) \approx 0$)
- ...but fails to generalize to new examples

Regularization

- A method for automatically controlling the complexity of the learned hypothesis
- Idea: penalize for large values of θ_j
 - Can incorporate into the cost function
 - Works well when we have a lot of features, each that contributes a bit to predicting the label
- Can also address overfitting by eliminating features (either manually or via model selection)

Regularization

Linear regression objective function

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

model fit to data regularization

- λ is the regularization parameter ($\lambda \geq 0$)
- No regularization on θ_0 !

Understanding Regularization

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

- Note that $\sum_{j=1}^{d} \theta_j^2 = \|\boldsymbol{\theta}_{1:d}\|_2^2$ - This is the magnitude of the feature coefficient vector!
- We can also think of this as: $\sum_{j=1}^{d} (\theta_j - 0)^2 = \|\boldsymbol{\theta}_{1:d} - \vec{\mathbf{0}}\|_2^2$
 - L₂ regularization pulls coefficients toward 0

Understanding Regularization

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

• What happens if we set λ to be huge (e.g., 10¹⁰)?



Based on example by Andrew Ng

Cost Function

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

- Fit by solving $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$
- Gradient update:

$$\frac{\partial}{\partial \theta_0} J(\theta) \quad \theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta} \left(\boldsymbol{x}^{(i)} \right) - y^{(i)} \right) \\
\frac{\partial}{\partial \theta_j} J(\theta) \quad \theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta} \left(\boldsymbol{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} - \lambda \theta_j$$

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

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta} \left(\boldsymbol{x}^{(i)} \right) - y^{(i)} \right)$$
$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta} \left(\boldsymbol{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} - \lambda \theta_j$$

• We can rewrite the gradient step as:

$$\theta_j \leftarrow \theta_j \left(1 - \alpha \lambda\right) - \alpha \frac{1}{n} \sum_{i=1}^n \left(h_{\theta}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right) x_j^{(i)}$$

• To incorporate regularization into the closed form solution:



• To incorporate regularization into the closed form solution:

$$\boldsymbol{\theta} = \left(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} + \lambda \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right)^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y}$$

- Can derive this the same way, by solving $\frac{\partial}{\partial \theta} J(\theta) = 0$
- Can prove that for λ > 0, inverse exists in the equation above