

# Linear Regression

Aarti Singh (Instructor), HMW-Alexander (Noter)

January 31, 2017

---

[Back to Index](#)

---

## Contents

<b>1 Discrete to Continuous Labels</b>	<b>1</b>
1.1 Task	1
1.2 Performance Measure	1
1.3 Bayes Optimal Rule	2
<b>2 Machine Learning Algorithm</b>	<b>2</b>
2.1 Empirical Risk Minimization (model-free)	2
<b>3 Linear Regression</b>	<b>2</b>
3.1 Gradient Descent	3
3.2 If $AA^T$ is not invertible	3
3.2.1 Regularized Least Squares	3

## Resources

- [Lecture](#)
- 

## 1 Discrete to Continuous Labels

From classification to regression

### 1.1 Task

Given  $X \in \mathcal{X}$ , predict  $Y \in \mathcal{Y}$ , Construct prediction rule  $f : \mathcal{X} \rightarrow \mathcal{Y}$

### 1.2 Performance Measure

- Quantifies knowledge gained.
- Measure of closeness between true label  $Y$  and prediction  $f(X)$ 
  - 0/1 loss:  $loss(Y, f(X)) = 1_{f(X) \neq Y}$ . Risk: probability of error
  - square loss:  $loss(Y, f(X)) = (f(X) - Y)^2$ . Risk: mean square error
- How well does the predictor perform on average?

$$\text{Risk } R(f) = \mathbb{E}[loss(Y, f(X))], (X, Y) \sim P_{XY}$$

### 1.3 Bayes Optimal Rule

- ideal goal: Construct prediction rule  $f^* : \mathcal{X} \rightarrow \mathcal{Y}$

$$f^* = \arg \min_f E_{XY}[\text{loss}(Y, f(X))]$$

(Bayes optimal rule)

- Best possible performance:

$$\forall f, R(f^*) \leq R(f)$$

(Bayes Risk)

Problem:  $P_{XY}$  is unknown.

Solution: Training data provides a glimpse of  $P_{XY}$

(observed)  $\{(X_i, Y_i)\} \sim_{i.i.d} P_{XY}$  unknown

## 2 Machine Learning Algorithm

- Model based approach: use data to learn a model for  $P_{XY}$
- Model-free approach: use data to learn mapping directly

### 2.1 Empirical Risk Minimization (model-free)

- Optimal predictor:

$$f^* = \arg \min_f \mathbb{E}[(f(X) - Y)^2]$$

- Empirical Minimizer:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X) - Y)^2$$

$\mathcal{F}$  is the class of predictors:

- Linear
- Polynomial
- Nonlinear

## 3 Linear Regression

$$f(\vec{X}) = \sum_{i=0}^p \beta_0 X^i = \vec{X}^T \vec{\beta}, \text{ where } X^0 = 1, \vec{\beta} = [\beta_0, \dots, \beta_p]^T$$

$$\hat{\vec{\beta}} = \arg \min_{\vec{\beta}} (A^T \vec{\beta} - \vec{Y})^T (A^T \vec{\beta} - \vec{Y}), \text{ where } A = [\vec{X}_1, \dots, \vec{X}_n]$$

$$J(\beta) = (A^T \vec{\beta} - \vec{Y})^T (A^T \vec{\beta} - \vec{Y})$$

$$\begin{aligned} \frac{\partial J(\vec{\beta})}{\partial \vec{\beta}} &= \frac{\partial (A^T \vec{\beta} - \vec{Y})^T (A^T \vec{\beta} - \vec{Y})}{\partial \vec{\beta}} \\ &= \frac{\partial (\vec{\beta}^T A A^T \vec{\beta} - \vec{\beta}^T A \vec{Y} - \vec{Y}^T A^T \vec{\beta} + \vec{Y}^T \vec{Y})}{\partial \vec{\beta}} \\ &= (A A^T + (A A^T)^T) \vec{\beta} - A \vec{Y} - A \vec{Y} \\ &= 2 A A^T \vec{\beta} - 2 A \vec{Y} = 0 \\ &\Rightarrow A A^T \vec{\beta} = A \vec{Y} \\ &\Rightarrow \hat{\vec{\beta}} = (A A^T)^{-1} A \vec{Y}, \text{ if } A A^T \text{ is invertible} \end{aligned}$$

### 3.1 Gradient Descent

Even when  $AA^T$  is invertible, might be computationally expensive if  $A$  is huge; however,  $J(\vec{\beta})$  is convex<sup>1</sup> in  $\beta$ .

Minimum of a convex function can be reached by gradient descent algorithm:

- Initialize: pick  $\vec{w}$  at random

- Gradient:

$$\nabla_{\vec{w}} l(\vec{w}) = \left[ \frac{\partial l(\vec{w})}{\partial w_0}, \dots, \frac{\partial l(\vec{w})}{\partial w_d} \right]^T$$

- Update rule:

$$\Delta \vec{w} = \eta \nabla_{\vec{w}} l(\vec{w})$$

,

$$w_i^{t+1} \leftarrow w_i^t - \eta \frac{\partial l(\vec{w})}{\partial w_i} \Big|_t$$

- Stop: when some criterion met  $\frac{\partial l(\vec{w})}{\partial w_i} \Big|_t < \epsilon$

### 3.2 If $AA^T$ is not invertible

$\text{Rank}(AA^T)$  = number of non-zero eigenvalues of  $AA^T$  = number of non-zero singular values of  $A \leq \min(n, p)$  since  $A$  is  $n \times p$

$$A = U\Sigma V^T \Rightarrow AA^T = U\Sigma^2 U^T \Rightarrow AA^T U = U\Sigma^2$$

#### 3.2.1 Regularized Least Squares

Ridge Regression (l2 penalty)

$$\begin{aligned} \hat{\vec{\beta}}_{MAP} &= \arg \min_{\vec{\beta}} (A^T \vec{\beta} - \vec{Y})^T (A^T \vec{\beta} - \vec{Y}) + \lambda \vec{\beta}^T \vec{\beta} \quad (\lambda \geq 0) \\ &= (AA^T + \lambda I)^{-1} A^T \vec{Y} \end{aligned} \quad (1)$$

$(AA^T + \lambda I)$  is invertible if  $\lambda > 0$ . Proof:

- the symmetric matrix  $AA^T$  is positive-semidefinite matrix, because a matrix is positive-semidefinite iff it arises as the Gram matrix of some set of vectors<sup>2</sup>.

- $\therefore \forall \lambda > 0$  and  $\vec{x} \neq \vec{0}$ ,

$$\begin{aligned} \vec{x}^T (AA^T) \vec{x} &= (A^T \vec{x})^T (A^T \vec{x}) \geq 0 \\ \vec{x}^T (AA^T + \lambda I) \vec{x} &= \vec{x}^T (AA^T) \vec{x} + \lambda \vec{x}^T \vec{x} > 0 \end{aligned}$$

- $\therefore (AA^T + \lambda I)$  is positive definite.
- $\therefore$  the eigenvalues of  $B = (AA^T + \lambda I)$  are all positive.

$$B\vec{v} = \lambda \vec{v} \Rightarrow \vec{v}^T B \vec{v} = \lambda > 0$$

- $\therefore (AA^T + \lambda I)$  is invertible if  $\lambda > 0$

<sup>1</sup>A function is called convex if the line joining any two points on the function does not go below the function on the interval formed by these two points.

<sup>2</sup>In contrast to the positive-definite case, these vectors need not be linearly independent.