

Regression

Dr. Fayyaz ul Amir Afsar Minhas

PIEAS Biomedical Informatics Research Lab Department of Computer and Information Sciences Pakistan Institute of Engineering & Applied Sciences PO Nilore, Islamabad, Pakistan http://faculty.pieas.edu.pk/fayyaz/

Regression

- Estimate the relationship among variables
 - Dependent variables
 - Independent variables
- Used in prediction and forecasting
- Mathematical formulation
 - Model: $Y = f(X; w) + \epsilon$
 - Linear

• $f(\mathbf{x}_i) = b + w_1 x_i^{(1)} + w_2 x_i^{(2)} + \dots + w_d x_i^{(d)}$

- Objective is to estimate the parameters w





Linear Regression

•
$$f(\mathbf{x}_i) = b + w_1 x_i^{(1)} + w_2 x_i^{(2)} + \dots + w_d x_i^{(d)}$$

This implies •

$$- f(\mathbf{x}_i) = [\mathbf{x}_i^T \quad 1]\mathbf{w}'$$

- Note that we do not have an explicit bias term anymore —
- For N points with

$$- y_1 = f(x_1) = [x_1^T \quad 1]w' - y_2 = f(x_2) = [x_2^T \quad 1]w'$$

$$X_{(N \times (d+1))} = \begin{bmatrix} x_1 \\ x_2^{(1)} \\ \vdots \\ x^{(1)} \end{bmatrix}$$

$$\begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(d)} & 1 \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(d)} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_N^{(1)} & x_N^{(2)} & \cdots & x_N^{(d)} & 1 \end{bmatrix}$$

(d)

 W_1

 $\boldsymbol{w'}_{((d+1)\times 1)} = \begin{vmatrix} w_2 \\ \vdots \\ w_d \end{vmatrix}$

$$- y_N = f(\mathbf{x}_N) = \begin{bmatrix} \mathbf{x}_N^T & 1b \end{bmatrix}$$

- In Matrix form •
 - y = Xw'•

$$\boldsymbol{y}_{(N\times 1)} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Linear Regression

• It can also be written as:

$$- Xw' = y$$

• If X is square (N=d+1) then the solution to the above equation is

$$-w'=X^{-1}y$$

• Example

$$-X = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix}, y = \begin{bmatrix} 3.5 \\ 4.75 \end{bmatrix}$$

- Thus: $w' = \begin{bmatrix} 1.25 \\ 2.25 \end{bmatrix}$





Linear Regression: Least Squares solution

- However, having a square X is very restrictive
- If X is not square, we can use a pseudo-inverse

$$Xw' = y$$
$$X^{T}Xw' = X^{T}y$$
$$w' = (X^{T}X)^{-1}X^{T}y$$
$$w' = X^{+}y$$

- $X^+ = (X^T X)^{-1} X^T$ is the pseudo-inverse
- Used to solve over-determined systems

 More constraints than parameters

Linear Regression: Least Squares solution

- Compute the pseudo-inverse and plot
- Note that the line isn't passing through all points



Properties of the least-squares solution

- The previous line represents a least squares (LS) solution to the regression problem
- It minimizes the mean square error of the prediction
- The mean square error resulting from a specific weight vector can be written as (ignoring bias):

$$- E(w) = \frac{1}{N} \sum_{i=1}^{N} (f(x_i) - y_i)^2 = \frac{1}{N} \sum_{i=1}^{N} (x_i^T w - y_i)^2$$

- Thus the LS learning problem can written as:

$$min_{w}\boldsymbol{E}(\boldsymbol{w}) = \frac{1}{N}\sum_{i=1}^{N}(x_{i}^{T}\boldsymbol{w}-y_{i})^{2}$$

- In Matrix form, $E(w) = \frac{1}{N} (Xw y)^T (Xw y)$
- If we differentiate E(w) with respect to w and substituting it to zero we get: Xw = y
 - We can now solve for **w** to get a closed form least square solution

Least squares visualized

- The thick red lines indicate the error corresponding to each data point
- Least square solution minimizes the sum of the square lengths of all red lines



Problems with least squares solution

- Due to squaring of the error of each data point the least square solution is very sensitive to outliers
- It gives only a linear solution



Support Vector Regression

- Hard Form
 - All errors must be within a user specified threshold
 - Minimize the norm of the weight vector





Penalty/loss function

Support Vector Regression

- Soft Form
 - Errors must be within a user specified threshold or penalize them linearly (instead of quadratically as in the LS solution)
 - Minimize the norm of the weight vector
- More robust!

$$\min_{w,b,\xi \ge 0} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

Such that for all i :
$$|(w^T x_i + b) - y_i| \le \epsilon + \xi_i$$

3+

-8



CIS 621: Machine Learning

So what is min $||w||^2$ doing here

 Geometric interpretation of margin maximization in regression



SVR: Dual Form

$$max_{\alpha} - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} (\alpha_{i}^{+} - \alpha_{i}^{-})(\alpha_{j}^{+} - \alpha_{j}^{-})\mathbf{x}_{i}^{T}\mathbf{x}_{j} - \epsilon \sum_{i=1}^{N} (\alpha_{i}^{+} + \alpha_{i}^{-}) - \sum_{i=1}^{N} y_{i}(\alpha_{i}^{+} - \alpha_{i}^{-})$$

Subject to:

$$0 \le \alpha_i^+ \le C, 0 \le \alpha_i^- \le C$$
$$\sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) = 0$$

 We can now apply the kernel trick and use it for non-linear regression

Kernels in SVR

 In the kernel space the SVR is fitting a line which corresponds to an arbitrary curve in the original feature space



SVR in action



Figure 13.7 The fitted regression line to data points shown as crosses and the ϵ -tube are shown ($C = 10, \epsilon = 0.25$). There are three cases: In (a), the instance is in the tube; in (b), the instance is on the boundary of the tube (circled instances); in (c), it is outside the tube with a positive slack, that is, $\xi_{+}^{t} > 0$ (squared instances). (b) and (c) are support vectors. In terms of the dual variable, in (a), $\alpha_{+}^{t} = 0, \alpha_{-}^{t} = 0$, in (b), $\alpha_{+}^{t} < C$, and in (c), $\alpha_{+}^{t} = C$.

SVR with quadratic kernel



Figure 13.8 The fitted regression line and the ϵ -tube using a quadratic kernel are shown ($C = 10, \epsilon = 0.25$). Circled instances are the support vectors on the margins, squared instances are support vectors which are outliers.

SVR with RBF kernel





Figure 13.9 The fitted regression line and the ϵ -tube using a Gaussian kernel with two different spreads are shown ($C = 10, \epsilon = 0.25$). Circled instances are the support vectors on the margins, and squared instances are support vectors that are outliers.

CIS 621: Machine Learning

Multi-output regression

- When we need to predict more than one variable as the output
 - The output variables may not be independent of each other

$$h: \Omega_{X_1} \times \ldots \times \Omega_{X_m} \longrightarrow \Omega_{Y_1} \times \ldots \times \Omega_{Y_d}$$
$$\mathbf{x} = (x_1, \ldots, x_m) \longmapsto \mathbf{y} = (y_1, \ldots, y_d),$$

- Solutions?
 - Apply a regression method for each variable
 - Shortcoming?
 - » Correlations are ignored
- Borchani, Hanen, Gherardo Varando, Concha Bielza, and Pedro Larrañaga. 2015. "A Survey on Multi-Output Regression." *Wiley Int. Rev. Data Min. and Knowl. Disc.* 5 (5): 216–33. doi:10.1002/widm.1157.

Required

- Reading:
 - Section 13.10 in Alpaydin 2010
- Problems
 - Discuss least squares regression and SVR in terms of structural risk minimization (SRM)
 - Understand how we achieved the dual form
 - Understand how we got from $|(w^T x_i + b) y_i| \le \epsilon + \xi_i$ to the following two constraints

$$y_i - (w^T x_i + b) \le \epsilon + \xi_i^+ (w^T x_i + b) - y_i \le \epsilon + \xi_i^-$$

We want to make a machine that will be proud of us.

- Danny Hillis

CIS 621: Machine Learning