

Some useful linear algebra & linear regression

Zhenzhen Kou

Thursday, January 20, 2005

Some useful linear algebra

- Scalar, Vector, Matrix
- Basic operations (+, -, *)
- Dot products (Inner product), Length of a vector: $\| \mathbf{x} \|$
- Transpose, Inverse
- Matrix calculus

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A}) = \mathbf{A}, \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{A} \mathbf{x}) = \mathbf{A}^T, \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}$$

The chain rule:
$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}},$$

The regression problem

- **Instances:** $\langle \mathbf{x}_j, t_j \rangle$
- **Learn:** Mapping from \mathbf{x} to $t(\mathbf{x})$
 - Find coeffs $\mathbf{w} = \{w_1, \dots, w_k\}$ in

$$\underbrace{t(\mathbf{x})}_{\text{data}} \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$$

- Precisely, minimize the **sum of squares residua (SSR)**:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

Why minimize SSR?

- Learned in class: maximizing log-likelihood in a Gaussian model

$$P(t \mid \mathbf{x}, \mathbf{w}, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[t - \sum_i w_i h_i(\mathbf{x})]^2}{2\sigma^2}}$$

- Minimizing SSR is also called **Least Square Error**

Linear regression

- Problem: Find coeffs $\mathbf{w} = \{w_1, \dots, w_k\}$ in

$$\underbrace{t(\mathbf{x})}_{\text{data}} \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$$

- Start from linear regression

$$\mathbf{y} = t(\mathbf{x}) = \sum_i w_i x_i = \mathbf{w}^T \mathbf{x}$$

- A trick to deal with the constant term:

$$\mathbf{x} = (x_1, \dots, x_{m-1}, 1)^T$$

Linear regression

Dataset has form

$$\begin{array}{cc} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \\ \vdots & \vdots \\ \cdot & \cdot \\ x_R & y_R \end{array}$$

Write matrix X and Y thus:

$$\mathbf{X} = \begin{bmatrix} \dots \mathbf{x}_1 \dots \\ \dots \mathbf{x}_2 \dots \\ \vdots \\ \dots \mathbf{x}_R \dots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ & & \vdots & \\ x_{R1} & x_{R2} & \dots & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_R \end{bmatrix}$$

Each row of X represents a record

Let \mathbf{Y} represent the output data and \mathbf{X} represent the input data, the linear regression model assumes a vector \mathbf{w} such that

$$\mathbf{Y} = \mathbf{X} \cdot \mathbf{w} \quad \text{Here } \mathbf{w} = \begin{bmatrix} w_1 \\ w_1 \\ \dots \\ w_m \end{bmatrix}$$

Linear regression

- To find the maximum likelihood estimation of \mathbf{w} is to look for \mathbf{w} which minimizes

$$\|\mathbf{X}\mathbf{w} - \mathbf{Y}\|^2 = (\mathbf{X}\mathbf{w} - \mathbf{Y})^T (\mathbf{X}\mathbf{w} - \mathbf{Y})$$

a concave function of \mathbf{w}

- Solve
$$\frac{\partial (\mathbf{X}\mathbf{w} - \mathbf{Y})^T (\mathbf{X}\mathbf{w} - \mathbf{Y})}{\partial \mathbf{w}} = \mathbf{0}$$

Apply the chain rule, we get:
$$\frac{\partial (\mathbf{X}\mathbf{w} - \mathbf{Y})^T}{\partial \mathbf{w}} \cdot 2(\mathbf{X}\mathbf{w} - \mathbf{Y}) = \mathbf{0}$$

i.e.
$$\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{Y}) = \mathbf{0}$$

- Therefore
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

What if $\mathbf{X}^T \mathbf{X}$ does not have inverse?
Pseudoinverse!

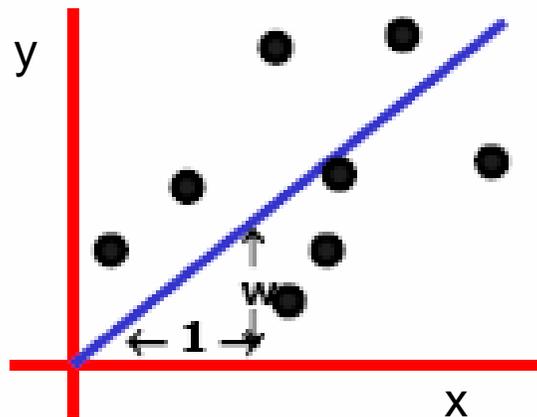
Linear regression

The max. likelihood w is $w = (X^T X)^{-1} (X^T Y)$

$X^T X$ is an $m \times m$ matrix: i, j 'th elt is $\sum_{k=1}^R x_{ki} x_{kj}$

$X^T Y$ is an m -element vector: i 'th elt $\sum_{k=1}^R x_{ki} y_k$

Linear Regression



DATASET

| inputs | outputs |
|-------------|-------------|
| $x_1 = 1$ | $y_1 = 1$ |
| $x_2 = 3$ | $y_2 = 2.2$ |
| $x_3 = 2$ | $y_3 = 2$ |
| $x_4 = 1.5$ | $y_4 = 1.9$ |
| $x_5 = 4$ | $y_5 = 3.1$ |

The regression problem

What if non-linear

- **Learn:** Mapping from \mathbf{x} to $t(\mathbf{x})$
 - Find coeffs $\mathbf{w}=\{w_1, \dots, w_k\}$ in

$$\underbrace{t(\mathbf{x})}_{\text{data}} \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$$

- What if $h(\mathbf{x})$ is a non-linear function?
 - Logistic regression

What you need to know

- Linear regression
 - Optimizing sum squared error
 - A linear regression has a close form solution
- More about regression
 - Andrew Moore's tutorials at <http://www-2.cs.cmu.edu/~awm/tutorials/>

Questions?

