

Linear Algorithms: An Introduction

Mia Feng

March 22, 2019

- Definition
- Classification
 - Perceptron
 - Logistic Regression
 - SVM
- Regression
 - OLS
 - Bayes Regression
- Comparision

What are LAs?

Linear Algorithms

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b. \quad (1)$$

To be clear,

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \quad (2)$$

where $\mathbf{w} = (w_1; w_2; \dots; w_d)$, $\mathbf{x} = (x_1; x_2; \dots; x_d)$.

Tasks

- classification
- regression

- Set hypothesis
e.g. $y = x_1 + 2x_2$
- Define cost function
e.g. RMSE, MSE
Accurate, Recall
- Find the optimal solution
e.g. gradient descent (GD) \Rightarrow Large dataset?
stochastic gradient descent (SGD) \Rightarrow optimal?
batch gradient descent \Rightarrow reasonable?
- Evaluate and model selection: cross validation + regularization
hold-out \Rightarrow accidental error?
leave-one-out \Rightarrow large dataset?
k-fold \Rightarrow reasonable?

Symbols and Notation

- J : cost function/control function
- $f(x)$: the hypothesis, a function of x with parameters w , where w is a scalar or a vector.
- b : intercept (scalar)
- x : feature, which is a vector or a scalar
- X : collection of features.
- y : label, which is a vector or a scalar
- ε : Gaussian noise, aka follows a Gaussian distribution with zero mean and variance σ^2
- η : learning rate

Perceptron-1

Idea: Minimize the number of misclassified samples.

$$J(\mathbf{w}, b) = - \sum_{\mathbf{x}_i \in X} y_i (\mathbf{w}^T \mathbf{x}_i + b). \quad (3)$$

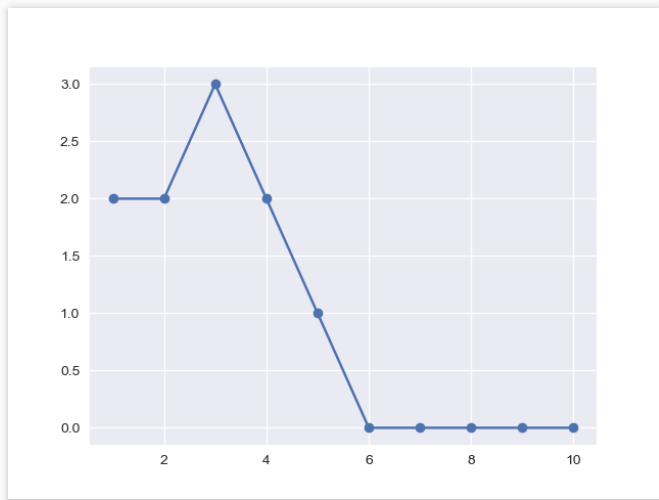
Solve: SGD

1. Initialize \mathbf{w}, b .
2. Repeat until no misclassified samples{
 - a. Choose (\mathbf{x}_i, y_i) randomly,
 - b. If $y_i (\mathbf{w}^T \mathbf{x}_i + b) \leq 0$:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i, \\ b &\leftarrow b + \eta y_i. \end{aligned} \quad (4)$$

}

Perceptron-2



The number of misclassified samples in each iteration.

Idea: Log MLE

$$J(\mathbf{w}) = - \sum_{i=1}^m y_i \ln f(\mathbf{x}_i) + (1 - y_i) \ln (1 - f(\mathbf{x}_i)). \quad (5)$$

Solve: GD

1. Initialize \mathbf{w}, b .
2. Repeat until convergence{

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \eta \frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i) \mathbf{x}_i, \\ b &\leftarrow b - \eta \frac{1}{m} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i). \end{aligned} \quad (6)$$

}

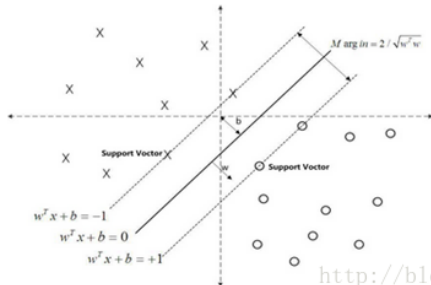
Idea: Maximize the geometry distance between hyperplanes

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2, \\ \text{s.t.} \quad & y_i(\mathbf{w}\mathbf{x}_i + b) - 1 \geq 0, \quad i = 1, 2, \dots, m; \end{aligned} \quad (7)$$

to be simpler,

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^m \alpha_i, \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0, \\ & \alpha_i \geq 0, \quad i = 1, 2, \dots, m. \end{aligned} \quad (8)$$

如下图所示，中间的实线便是寻找到的最优超平面（Optimal Hyper Plane），其到两条虚线的距离相等，这个距离便是几何间隔 $\tilde{\gamma}$ ，两条虚线之间的距离等于 $2\tilde{\gamma}$ ，而虚线上的点则是支持向量。由于这些支持向量刚好在边界上，所以它们满足 $y(w^T x + b) = 1$ （还记得我们把 functional margin 定为 1 了吗？上节中：处于方便推导和优化的目的，我们可以令 $\hat{\gamma} = 1$ ），而对于所有不是支持向量的点，则显然有 $y(w^T x + b) > 1$ 。



<http://blog.csdn.net/puqutogether>

Hyperplane and support vector

SVM-3: See equations (p125-p131) in Book by Hang Li.

算法 7.5 (SMO 算法)

输入: 训练数据集 $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, 其中, $x_i \in \mathcal{X} = \mathbf{R}^n$, $y_i \in \mathcal{Y} = \{-1, +1\}$, $i = 1, 2, \dots, N$, 精度 ε ;

输出: 近似解 $\hat{\alpha}$.

(1) 取初值 $\alpha^{(0)} = 0$, 令 $k = 0$;

(2) 选取优化变量 $\alpha_1^{(k)}, \alpha_2^{(k)}$, 解析求解两个变量的最优化问题 (7.101) ~ (7.103), 求得最优解 $\alpha_1^{(k+1)}, \alpha_2^{(k+1)}$, 更新 α 为 $\alpha^{(k+1)}$;

(3) 若在精度 ε 范围内满足停机条件

$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, N$$

$$y_i \cdot g(x_i) = \begin{cases} \geq 1, & \{x_i \mid \alpha_i = 0\} \\ = 1, & \{x_i \mid 0 < \alpha_i < C\} \\ \leq 1, & \{x_i \mid \alpha_i = C\} \end{cases}$$

其中,

$$g(x_i) = \sum_{j=1}^N \alpha_j y_j K(x_j, x_i) + b$$

则转 (4); 否则令 $k = k + 1$, 转 (2);

(4) 取 $\hat{\alpha} = \alpha^{(k+1)}$.

SMO algorithm

Ordinary Least Square

Idea: Minimize MSE

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2. \quad (9)$$

Solve: GD

Repeat until convergence{

$$w_j = w_j + \eta \sum_{i=1}^m (y_i - f(\mathbf{x}_i)). \quad (10)$$

}

Idea: MAP

$$\log p(y|X) = -\frac{1}{2}y^T K^{-1}y - \frac{1}{2} \log |K| - \frac{m}{2} \log 2\pi, \quad (11)$$

where K is the covariance matrix of training data X .

For more details, read "GPR for ML" by Carl Edward Rasmussen (p8-p22).

Solve

Repeat until convergence or out of the iteration limitations.

cost function

- Regression: RMSE, MSE etc.
- Classification: the number of misclassified samples, recall, likelihood etc.

General

Derivative

- Determined:
 - Geometry: SVM
 - The number of misclassified samples: Perceptron
 - Euclidean distance: OLS
- Stochastic:
 - Frequency: Logistic, OLS(\leftrightarrow MLE if data follow unbiased Gaussian distribution)?
 - Bayesian: Bayesian regression.

Relation

Bayesian regression \leftrightarrow OLS + random noise
prior \leftrightarrow regularization