

ML Note

qhy

2017 年 7 月 31 日

目录

1	SVM	2
1.1	线性SVM	3
1.2	线性SVM2	3
1.3	Kernel SVM	4
1.4	SVM的应用与变种	4
2	Structure Learning	4
2.1	Structured Liner Model	5
2.2	structure SVM	7
2.2.1	Multi-class SVM	8

1 SVM

定义标签: $\hat{y} = \{+1, -1\}$, +1表示正类, -1表示负类定义输出:

$$g(x) = \begin{cases} +1 & , f(x) > 0 \\ -1 & , f(x) < 0 \end{cases} \quad (1)$$

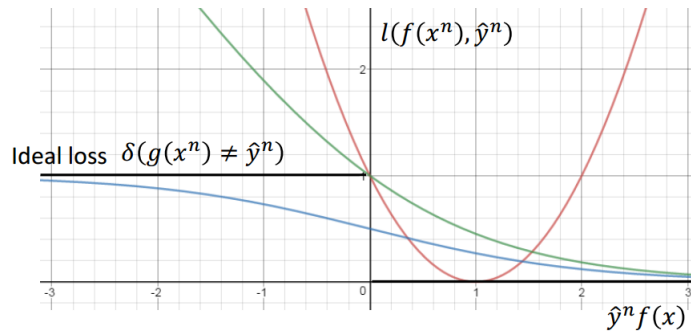
定义loss function:

$$L(f) = \sum_n I(g(x^n) \neq \hat{y}^n) \quad (2)$$

也就是统计分类错误的个数,从更一般的角度来看,就是每分类错误一个,就给一定的惩罚

分析: $\hat{y}^n f(x^n) > 0$ 的时候,分类正确,无论是正类还是负类,

如果把这个看做一个新的变量 t ,那么损失函数(一个数据的损失函数)就可以定义成 $t > 0$ 时无惩罚, $t < 0$ 的时候给出一定的惩罚。



理想的情况就是图1中的黑线,分类错误就给出惩罚。

由于此loss function使用梯度下降算法麻烦,故使用一个惩罚函数 l 代替 I :

$$L(f) = \sum_n I(f(x^n), \hat{y}^n) \quad (3)$$

红色线为**Square loss**:

$$l(f(x^n), \hat{y}^n) = (\hat{y}^n f(x^n) - 1)^2 \quad (4)$$

显然是不合理的,因为它对于分类正确的结果也给出了惩罚。

蓝色线为**sigmoid + square loss**:

$$l(f(x^n), \hat{y}^n) = (\sigma(\hat{y}^n f(x^n)) - 1)^2 \quad (5)$$

这样,得到的结果是合理的,对于分类正确,确定性高的给出低惩罚,对于分类错误,甚至错误的离谱的,则给出大的惩罚

但是有一个问题就是,对于分类错误的结果,它的惩罚很低。

绿色线为**Sigmoid + cross entropy**

$$l(f(x^n), \hat{y}^n) = \ln(1 + \exp(-\hat{y}^n f(x^n))) \quad (6)$$

蓝色线为**hinge loss**:

$$l(f(x^n), \hat{y}^n) = \max(0, 1 - \hat{y}^n f(x^n)) \quad (7)$$

对于分类正确的结果,只要 $\hat{y}^n f(x^n) < 1$ 的数据还是会给出惩罚,而SVM采用的就是HingeLoss从 $f(x)$ 来看, $f(x)$ 为分界线,只要在 $f(x)$ 上下距离为1的样本都会给出惩罚。

1.1 线性SVM

模型如下:

- Step 1: Function (Model)

$$f(x) = \sum_i w_i x_i + b = \begin{matrix} \text{New } w \\ \begin{bmatrix} w \\ b \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{New } x \\ \begin{bmatrix} x \\ 1 \end{bmatrix} \end{matrix} = w^T x$$

- Step 2: Loss function

$$L(f) = \sum_n l(f(x^n), \hat{y}^n) + \lambda \|w\|_2$$

convex (pointing to the loss function graph)

$$l(f(x^n), \hat{y}^n) = \max(0, 1 - \hat{y}^n f(x^n))$$

梯度下降:

$$\frac{\partial L(f)}{\partial w_i} = \sum_n \frac{-\delta(\hat{y}^n f(x^n) < 1) \hat{y}^n x_i}{c^n(w)} \quad w_i \leftarrow w_i - \eta \sum_n c^n(w) x_i^n$$

1.2 线性SVM2

至此,线性SVM结束,但是这个SVM和之前见到的SVM貌似不太一样?对偶问题,SMO,等等这些东东呢?

另一种描述:

Minimizing loss function L:

$$L(f) = \sum_n \varepsilon^n + \lambda \|w\|_2$$

$$\varepsilon^n = \max(0, 1 - \hat{y}^n f(x^n))$$

||

ε^n : slack variable
Quadratic programming problem

$$\varepsilon^n \geq 0$$

$$\varepsilon^n \geq 1 - \hat{y}^n f(x^n) \Rightarrow \hat{y}^n f(x^n) \geq 1 - \varepsilon^n$$

经过这样处理,原问题从最小化:

$$L(f) = \sum_n l(f(x^n), \hat{y}^n) + \lambda \|w\|_2$$

变成了最小化:

$$L(f) = \sum_n \varepsilon^n + \lambda \|w\|_2 \quad \text{s.t.} \quad \hat{y}^n f(x^n) \geq 1 - \varepsilon^n$$

进一步处理(去掉松弛变量),问题最小化:

$$L(f) = \lambda \|w\|_2 \quad \text{s.t.} \quad \hat{y}^n f(x^n) \geq 1$$

1.3 Kernel SVM

简单地从梯度下降的角度来看,那么最终得到的 w^* 一定是所有 x_i 的线性组合,i.e.

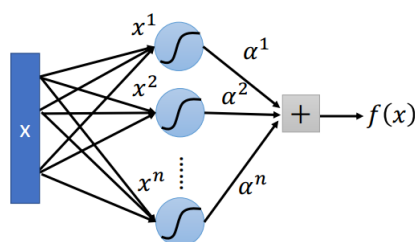
$$w^* = \sum_i a_i^* x_i \quad (8)$$

而由于使用**hinge loss**的原因,大部分的 $a_i^* = 0$,对于 $a_i^* \neq 0$ 的 x_i ,则称为**Support Vector**

那么最终的解的形式变成了 $f(x) = \sum_i a_i (x_i \cdot x)$

若引入核函数,则是 $f(x) = \sum_i a_i \cdot K(x_i \cdot x)$

而这个解的形式,则可以看成一个单隐层神经网络:



其中,每一个**Support Vector**表示一个神经元。

这里为了引入核函数,把问题变成了这种形式,参数变成了 a ,但一样可以用,梯度下降去解决,从这里来看,应该是由 a 决定哪些是 **支持向量**,而不是由支持向量,决定哪些 $a = 0$

为什么说RBF是对应的是无穷维的点积?

因为RBF使用泰勒展开之后,有无穷项,每一项代表若干个维度的分量分别相乘之和的结果。

怎么设计核函数? 满足 mercer 条件的任一函数都可以是核函数。

1.4 SVM的应用与变种

- SVR
- Ranking SVM
- one-class SVM
- Recursive Kernel

2 Structure Learning

之前我们接触的大部分都是输入是一个向量,输出是一个实数但如果我们要输出一个序列则怎么处理?

这里有一个大致的框架:

设 X 是输入, Y 是输出, $F(X, Y)$ 描述 X, Y 的相似度,或者是 X 输出 Y 的概率

对于训练过程,我们要做的就是最大化训练样本输入和输出的相似度,i.e.

$$\max F(X, Y) \quad (9)$$

对于测试过程,我们的输出就是在 Y 的空间中,与测试输入 $X, F(X, Y)$ 最大的那个 Y

那么,就有3个问题需要解决:

1. $F(X, Y)$ 是什么样子?
2. 怎么有效快速地输出 y ?
3. 怎么训练?

总的来说,就是把 X 和 Y 两个不同的东西,表示在同一个空间内,这样才能建立起度量,进而才有优化的方向。

这三个问题和 HMM 中的三个问题是一致的,因为HMM就是Structure Learning中的一种

2.1 Structured Liner Model

$F(x, y)$ 的形式: y 可以看成是 x 在特征空间中的投影,我们要做的就是在特征空间中的点距离 y 尽可能地近,或者是相似度尽可能大。

$$F(x, y) = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w \end{bmatrix} \cdot \begin{bmatrix} \phi_1(x, y) \\ \phi_2(x, y) \\ \phi_3(x, y) \\ \vdots \\ \phi(x, y) \end{bmatrix}$$

↓

$$F(x, y) = w \cdot \phi(x, y)$$

我们可以简单地定义损失函数就是 $F(x, \hat{y}) - F(x, y)$, 也就是当前预测的点到目标点的距离。

其中 \hat{y} 是 $F(x, y)$ 最大的 y ,我们只需要考虑最大的 y 即可,而不需要考虑其他。

至此,我们就可以用梯度下降的方法进行训练:

- **Input:** training data set $\{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^r, \hat{y}^r), \dots\}$
- **Output:** weight vector w
- **Algorithm:** Initialize $w = 0$
 - do
 - For each pair of training example (x^r, \hat{y}^r)
 - Find the label \tilde{y}^r maximizing $w \cdot \phi(x^r, y)$

$$\tilde{y}^r = \arg \max_{y \in Y} w \cdot \phi(x^r, y) \text{ (question 2)}$$
 - If $\tilde{y}^r \neq \hat{y}^r$, update w

$$w \rightarrow w + \phi(x^r, \hat{y}^r) - \phi(x^r, \tilde{y}^r)$$
 - until w is not updated ➡ We are done!

收敛性证明:对于线性可分的数据,我们可以证明,最多需要 $(\frac{R}{\delta})^2$ 次迭代就可以收敛。

其中, R 表示特征空间中任意两个 y 之间的最大距离, δ 表示margin

也就是说,迭代的次数和样本的规模无关。

$$w^k = w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n) \text{ (the relation of } w^k \text{ and } w^{k-1})$$

Proof that: The angle ρ_k between \hat{w} and w_k is smaller as k increases

i.e. Analysis $\cos \rho_k = \frac{\hat{w} \cdot w^k}{\|\hat{w}\| \cdot \|w^k\|}$

$$\begin{aligned} \hat{w} \cdot w^k &= \hat{w} \cdot (w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)) \\ &= \hat{w} \cdot w^{k-1} + \underbrace{\hat{w} \cdot \phi(x^n, \hat{y}^n)}_{\geq \delta} - \hat{w} \cdot \phi(x^n, \tilde{y}^n) \geq \hat{w} \cdot w^{k-1} + \delta \\ &\geq \delta \text{ (Separable)} \end{aligned}$$

即,我们要证明,随着迭代的过程, w 是不断向最优解靠近的,也就是夹角不断变小

$$\hat{w} \cdot w^k \geq \hat{w} \cdot w^{k-1} + \delta \longrightarrow \hat{w} \cdot w^k \geq k\delta$$

$$\begin{aligned} \|w^k\|^2 &= \|w^{k-1} + \phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)\|^2 \\ &= \|w^{k-1}\|^2 + \underbrace{\|\phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n)\|^2}_{>0} + \underbrace{2w^{k-1} \cdot (\phi(x^n, \hat{y}^n) - \phi(x^n, \tilde{y}^n))}_{? < 0 \text{ (mistake)}} \end{aligned}$$

Assume the distance between any two feature vector is smaller than R

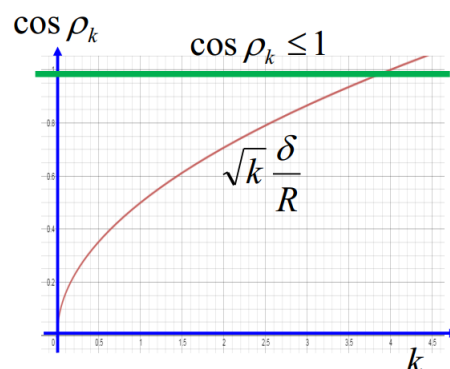
$$\begin{aligned} \|w^1\|^2 &\leq \|w^0\|^2 + R^2 = R^2 \\ \|w^2\|^2 &\leq \|w^1\|^2 + R^2 \leq 2R^2 \\ &\dots \\ \|w^k\|^2 &\leq kR^2 \end{aligned}$$

$$\cos \rho_k = \frac{\hat{w} \cdot w^k}{\|\hat{w}\| \cdot \|w^k\|} \quad \hat{w} \cdot w^k \geq k\delta \quad \|w^k\|^2 \leq kR^2$$

$$\geq \frac{k\delta}{\sqrt{kR^2}} = \sqrt{k} \frac{\delta}{R}$$

$$\sqrt{k} \frac{\delta}{R} \leq 1$$

$$k \leq \left(\frac{R}{\delta}\right)^2$$



这里我们并没有证明 $\cos \rho$ 是增大到1的,我们证明的是它的一个下界是不断增大到1,那么根据夹逼定理得证。

值得注意的是,这里有一个参数 δ ,这个需要事先人为设定,当约束比较严格的是,则迭代次数较多,当约束比较松的时候,则迭代次数少。

δ :对于这个参数的理解,我们并不要求整个算法的输入 X 在特征空间中的表示一定要和 Y 的表示重合,只要它们的差值在一定的范围内都是可接受的,因此,这个参数也叫margin

实际上,上述的模型中loss function并没有体现出来,新的loss function如下:

$$C^n = \max_y [w \cdot \phi(x^n, y)] - w \cdot \phi(x^n, \hat{y}^n)$$

$$C^n = \max_y [\Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)] - w \cdot \phi(x^n, \hat{y}^n)$$

2.2 structure SVM

光从定义上来看,SVM的框架不变,只是SVM衡量距离的部分变成了 $F(x, y)$

现在,我们根据上面的loss function转换到我们熟悉的SVM框架上来:

Find w minimizing C

$$C = \frac{1}{2} \|w\|^2 + \lambda \sum_{n=1}^N C^n$$

$$C^n = \max_y [\Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)] - w \cdot \phi(x^n, \hat{y}^n)$$

Are they equivalent? We want to minimize C

For $\forall y$:

$$C^n + w \cdot \phi(x^n, \hat{y}^n) = \max_y [\Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)]$$

$$C^n + w \cdot \phi(x^n, \hat{y}^n) \geq \Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)$$

$$w \cdot \phi(x^n, \hat{y}^n) - w \cdot \phi(x^n, y) \geq \Delta(\hat{y}^n, y) - C^n$$

Find w minimizing C

$$C = \frac{1}{2} \|w\|^2 + \lambda \sum_{n=1}^N C^n$$

$$C^n = \max_y [\Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)] - w \cdot \phi(x^n, \hat{y}^n)$$

Find $w, \varepsilon^1, \dots, \varepsilon^N$ minimizing C

$$C = \frac{1}{2} \|w\|^2 + \lambda \sum_{n=1}^N \varepsilon^n$$

For $\forall n$:

For $\forall y$:

$$w \cdot \phi(x^n, \hat{y}^n) - w \cdot \phi(x^n, y) \geq \Delta(\hat{y}^n, y) - \varepsilon^n$$

Slack variable

注意:如果去掉了minimizing,两个问题两个问题是不等价的,前者可以得到唯一的确定的解,而后者会得到一个范围。

好了,到这里我们就可以使用SVM中的套路来解决问题了,但是有一个问题,那就是当数据规模过大的时候,约束条件特别多,二次规划问题的代价就特别大。

这里给出了Cutting Plane Algorithm

Given training data: $\{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^N, \hat{y}^N)\}$

Working Set $\mathbb{A}^1 \leftarrow \text{null}, \mathbb{A}^2 \leftarrow \text{null}, \dots, \mathbb{A}^N \leftarrow \text{null}$

Repeat

$w \leftarrow \text{Solve a QP with Working Set } \mathbb{A}^1, \mathbb{A}^2, \dots, \mathbb{A}^N$

For each training data (x^n, \hat{y}^n) :

$$\bar{y}^n = \arg \max_y [\Delta(\hat{y}^n, y) + w \cdot \phi(x^n, y)]$$

find the most violated constraints

Update working set $\mathbb{A}^n \leftarrow \mathbb{A}^n \cup \{\bar{y}^n\}$

Until $\mathbb{A}^1, \mathbb{A}^2, \dots, \mathbb{A}^N$ doesn't change any more

Return w

这个算法的思想是,并不是所有的约束条件都是有用的,实际上对解的约束只有少部分条件是直接作用的,其他条件都是在满足直接条件下恰好满足的。

因此这个算法就是不断地迭代,随着迭代不断把约束条件加入到“直接”约束空间(初始为无约束)中,知道当求出来的解能满足所有约束条件位置。

这里每次迭代都加入约束最严格的那个条件,那么怎么定义这个最严格?这里把距离直接取 $F(x, y) - F(x, \hat{y})$ 最大的 y 作为最严格的约束条件

2.2.1 Multi-class SVM

多分类的SVM也属于 *StructureSVM* 的一种,输出的一个维度表示是否是那一类。

而且多分类的问题可以简化:在求损失的时候,只要求对应类标签和对应类参数下预测值的差距即可。也就是第 i 类的训练样本,只对第 i 类的参数有影响