

7. Soft-Margin SVM & Kernels

Montag, 3. Dezember 2018 17:15

Problem 1: What is the connection between soft-margin SVM and logistic regression?

Soft-margin SVM:

$$\underset{\vec{w}, b, \xi}{\text{minimize}} \quad f_0(\vec{w}, b, \xi) = \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \xi_i$$

$$\text{subject to} \quad y_i (\vec{w}^T \vec{x}_i + b) - 1 + \xi_i \geq 0 \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

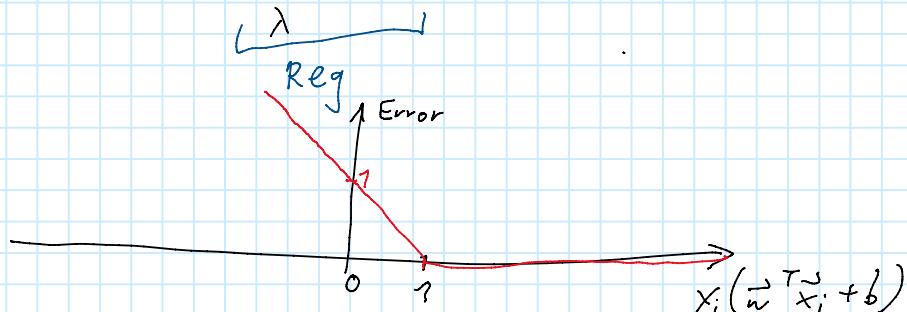
$$y_i (\vec{w}^T \vec{x}_i + b) - 1 + \xi_i = 0$$

$$\xi_i = \begin{cases} 1 - y_i (\vec{w}^T \vec{x}_i + b) & \text{if } y_i (\vec{w}^T \vec{x}_i + b) < 1 \\ 0 & \text{else} \end{cases}$$

rewrite as unconstrained opt. problem:

$$\underset{\vec{w}, b}{\text{minimize}} \quad \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \max(0, 1 - y_i (\vec{w}^T \vec{x}_i + b))$$

$$E(\vec{w}, b, C) = \underbrace{\frac{1}{2} C \vec{w}^T \vec{w}}_{\lambda} + \sum_i \text{Hinge}(y_i (\vec{w}^T \vec{x}_i + b))$$



Logistic regression:

$$p(y_i = 1 | \vec{x}, \vec{w}) = \sigma(\vec{w}^T \vec{x}_i + b)$$

$$p(y_i = -1 | \vec{x}, \vec{w}) = \sigma(-(\vec{w}^T \vec{x}_i + b)) = 1 - \sigma(\vec{w}^T \vec{x}_i + b)$$

$$\left[y_i \in \{0, 1\}: p(\vec{y} | \vec{X}, \vec{w}) = \prod_{i=1}^N \sigma(\vec{w}^T \vec{x}_i + b)^{y_i} \cdot (1 - \sigma(\vec{w}^T \vec{x}_i + b))^{1-y_i} \right]$$

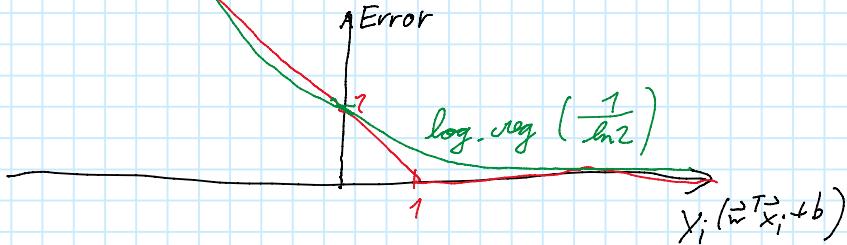
$$y_i \in \{-1, 1\} \quad p(\vec{y} | \vec{X}, \vec{w}) = \prod_{i=1}^N \underbrace{\sigma(y_i (\vec{w}^T \vec{x}_i + b))}_{\frac{1}{1+e^{-z}}}$$

$$E(\vec{w}, b) = -\ln(p(\vec{y} | \vec{X}, \vec{w})) = -\sum_{i=1}^N \ln((1 + e^{-y_i (\vec{w}^T \vec{x}_i + b)})^{-1}) =$$

$$E(\vec{w}, b) = -\ln(p(\vec{y} | X, \vec{w})) = -\sum_{i=1}^N \ln((1 + e^{-y_i(\vec{w}^T \vec{x}_i + b)})^{-1}) =$$

$$= \sum_{i=1}^N \ln(1 + e^{-y_i(\vec{w}^T \vec{x}_i + b)})$$

+ L2 Regularization: $E(\vec{w}, b, \lambda) = \lambda \vec{w}^T \vec{w} + E(\vec{w}, b)$



2 Gaussian kernel

Problem 2: One of the nice things about kernels is that new kernels can be constructed out of already given ones. Use the five kernel construction rules from the lecture to prove that the function

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|^2}{2\sigma^2}\right)$$

is a kernel.

$$\exp(K(\vec{x}, \vec{y})) = 1 + \sum_{n=1}^{\infty} \frac{1}{n!} K(\vec{x}, \vec{y})^n$$

$K(\vec{x}, \vec{y})^n$ Product ✓

$\frac{1}{n!} K(\vec{x}, \vec{y})^n$ Scalar mult ✓

$\sum_{n=1}^{\infty} \frac{1}{n!} K(\vec{x}, \vec{y})^n$ Sum ✓

$$K_3(\phi(\vec{x}), \phi(\vec{y})) \quad \phi(\vec{z}) = 1 \quad \checkmark$$

$$1 + \sum_{n=1}^{\infty} \frac{1}{n!} K(\vec{x}, \vec{y}) \quad \text{Sum} \quad \checkmark$$

$$\exp\left(-\frac{|\vec{x} - \vec{y}|^2}{2\sigma^2}\right) = \exp\left(-\frac{1}{2\sigma^2} (\vec{x} - \vec{y})^T (\vec{x} - \vec{y})\right) = \exp\left(-\frac{1}{2\sigma^2} (\vec{x}^T \vec{x} - 2 \vec{x}^T \vec{y} + \vec{y}^T \vec{y})\right) =$$

$$= \underbrace{\exp\left(-\frac{\vec{x}^T \vec{x}}{2\sigma^2}\right)}_{\phi(\vec{x})} \underbrace{\exp\left(-\frac{\vec{y}^T \vec{y}}{2\sigma^2}\right)}_{\phi(\vec{y})} \underbrace{\exp\left(\frac{-2\vec{x}^T \vec{y}}{\sigma^2}\right)}_{\frac{1}{\sigma^2} K(\vec{x}, \vec{y})}$$

$$\phi(\vec{z}) = \exp\left(-\frac{\vec{z}^T \vec{z}}{2\sigma^2}\right)$$

$$\frac{1}{\sigma^2} K(\vec{x}, \vec{y})$$

$$K_r(\vec{x}, \vec{y}) = \exp\left(\frac{\vec{x}^T \vec{y}}{\sigma^2}\right)$$

$$\rightarrow K_\ell(\vec{x}, \vec{y}) = \phi(\vec{x}) \cdot \phi(\vec{y})$$

Multiplication $K_\ell(\vec{x}, \vec{y}) \cdot K_r(\vec{x}, \vec{y}) \quad \checkmark$

3 Stacking feature maps

Suppose you have found a feature map $\theta : \mathbb{R}^n \rightarrow \mathbb{R}^m$ that transforms your data into a feature space in which a SVM with a Gaussian kernel works well. However computing the feature map $\theta(\mathbf{x})$ is computationally expensive and luckily you discover an efficient method to compute the scalar product $K(\mathbf{x}, \mathbf{y}) = \theta(\mathbf{x})^T \theta(\mathbf{y})$ in your feature space without having to compute $\theta(\mathbf{x})$ and $\theta(\mathbf{y})$ explicitly.

Problem 3: Show how you can use the scalar product $K(\mathbf{x}, \mathbf{y})$ to efficiently compute the Gaussian kernel in your feature space, that is

$$K_g(\theta(\mathbf{x}), \theta(\mathbf{y}))$$

where

$$K_g(\mathbf{a}, \mathbf{b}) = \exp\left(-\frac{|\mathbf{a} - \mathbf{b}|^2}{2\sigma^2}\right)$$

is the Gaussian kernel.

$$\begin{aligned} K_g(\theta(\tilde{\mathbf{x}}), \theta(\tilde{\mathbf{y}})) &= \exp\left(-\frac{|\theta(\tilde{\mathbf{x}}) - \theta(\tilde{\mathbf{y}})|^2}{2\sigma^2}\right) = \\ &= \exp\left(-\frac{(\theta(\tilde{\mathbf{x}}) - \theta(\tilde{\mathbf{y}}))^T (\theta(\tilde{\mathbf{x}}) - \theta(\tilde{\mathbf{y}}))}{2\sigma^2}\right) = \\ &= \exp\left(-\frac{\theta(\tilde{\mathbf{x}})^T \theta(\tilde{\mathbf{x}}) - 2\theta(\tilde{\mathbf{x}})^T \theta(\tilde{\mathbf{y}}) + \theta(\tilde{\mathbf{y}})^T \theta(\tilde{\mathbf{y}})}{2\sigma^2}\right) = \\ &= \exp\left(-\frac{K(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}) - 2K(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) + K(\tilde{\mathbf{y}}, \tilde{\mathbf{y}})}{2\sigma^2}\right) \\ &\Rightarrow \text{easy to compute} \end{aligned}$$

4 Unsuitable Kernels

Consider a SVM without slack variables using the kernel $K(x, y) = |x||y|$ where $x, y \in \mathbb{R}$.

Problem 4: Show that the function $K(x, y)$ is indeed a kernel and write down a feature space $\phi(x)$ corresponding to this kernel.

$$\phi(x) = |x| \quad K(x, y) = \phi(x)^T \phi(y) \quad \checkmark$$

Problem 5: Write down the set of classification functions that can be learned by an SVM using the kernel $K(x, y)$. Make sure that the set contains each classification function only once.

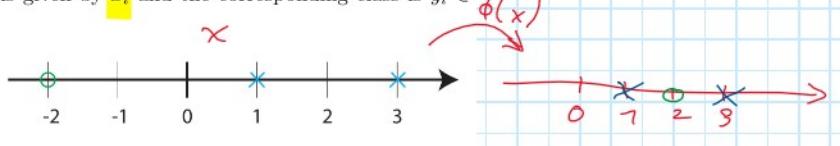
$$h_b(x) = \operatorname{sgn}(\phi(x) + b) = \operatorname{sgn}(|x| + b) = \begin{cases} 1 & \text{if } |x| + b < 0 \end{cases}$$

$$= \begin{cases} -1 & \text{if } |x| + b < 0 \\ 0 & \text{if } |x| + b = 0 \\ 1 & \text{if } |x| + b > 0 \end{cases}$$

$$S = \{h_b : b \in \mathcal{O}\} \cup \{h_1\}$$

Consider the following data set. The i th data point is given by x_i and the corresponding class is $y_i \in \{-1, +1\}$.

i	x_i	y_i
1	-2	-1
2	1	+1
3	3	+1



Problem 6: What would happen if you try to solve the dual problem for fitting an SVM without slack variables to this data set? Explain your answer.

$$\max \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j y_i y_j \alpha_i \alpha_j x_i x_j$$

$$\text{subject to } \sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

$$\boxed{\alpha_1 + \alpha_2 + \alpha_3} - \frac{1}{2} \left(\alpha_1^2 + 4\alpha_2^2 + 9\alpha_3^2 - 2 \cdot 2\alpha_1 \alpha_2 + 2 \cdot 3\alpha_1 \alpha_3 - 2 \cdot 3\alpha_2 \alpha_3 \right)$$

$\boxed{= 0} \text{ if } \alpha_1 + 3\alpha_3 = 2\alpha_2$

$$\sum_i \alpha_i y_i = \alpha_1 - \alpha_2 + \alpha_3 = 0 \Leftrightarrow \alpha_2 = \alpha_1 + \alpha_3$$

$$\frac{1}{2}\alpha_1 + \frac{3}{2}\alpha_3 = \alpha_1 + \alpha_3$$

$$\alpha_3 = \alpha_1$$

maximize with $\alpha_1 = \alpha_3 = \frac{\alpha_2}{2} \rightarrow \infty$

5 Kernelized k -nearest neighbors

To classify the point \mathbf{x} the k -nearest neighbors finds the k training samples $\mathcal{N} = \{\mathbf{x}^{(s_1)}, \mathbf{x}^{(s_2)}, \dots, \mathbf{x}^{(s_k)}\}$ that have the shortest distance $\|\mathbf{x} - \mathbf{x}^{(s_i)}\|_2$ to \mathbf{x} . Then the label that is mostly represented in the neighbor set \mathcal{N} is assigned to \mathbf{x} .

Problem 7: Formulate the k -nearest neighbors algorithm in feature space by introducing the feature map $\phi(\mathbf{x})$. Then rewrite the k -nearest neighbors algorithm so that it only depends on the scalar product in feature space $K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$.

$$\|\phi(\tilde{\mathbf{x}}) - \phi(\tilde{\mathbf{x}}^{(s_i)})\|_2^2 = (\phi(\tilde{\mathbf{x}}) - \phi(\tilde{\mathbf{x}}^{(s_i)}))^T (\phi(\tilde{\mathbf{x}}) - \phi(\tilde{\mathbf{x}}^{(s_i)})) =$$

$$= \phi(\tilde{\mathbf{x}})^T \phi(\tilde{\mathbf{x}}) - 2 \phi(\tilde{\mathbf{x}})^T \phi(\tilde{\mathbf{x}}^{(s_i)}) + \phi(\tilde{\mathbf{x}}^{(s_i)})^T \phi(\tilde{\mathbf{x}}^{(s_i)})$$

find $\tilde{\mathbf{x}}^{(s_i)}$ that minimize this. (k samples)

$$\boxed{-2K(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}^{(s_i)}) + K(\tilde{\mathbf{x}}^{(s_i)}, \tilde{\mathbf{x}}^{(s_i)})}$$