

1. Докажите, что задача логистической регрессии является выпуклой задачей машинного обучения (строго обоснуйте каждый переход, например, если говорите, что выпукла, так как композиция линейной функции и выпуклой, то обозначьте, какая функция линейна, а какая выпуклая)

Множество гипотез

$$H_{\text{sig}} = \phi_{\text{sig}} \circ L_d = \{\mathbf{x} \mapsto \phi_{\text{sig}}(\langle \mathbf{w}, \mathbf{x} \rangle) : \mathbf{w} \in \mathbb{R}^d\}.$$

будет являться выпуклым т.к. $H = \mathbb{R}^d$

Сама задача имеет вид

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)).$$

Доменом будет $Z = X \times Y = \mathbb{R}^{d+1}$

Рассмотрим ф-ию потерь на этом домене.

$$\ell(h_{\mathbf{w}}, (\mathbf{x}, y)) = \log(1 + \exp(-y \langle \mathbf{w}, \mathbf{x} \rangle)).$$

И покажем, что она выпукла.

Рассмотрим ф-ию $g(x) = \log(1 + \exp(x))$. И найдем её вторую производную.

$$g'(x) = \frac{\exp(x)}{1 + \exp(x)}$$

$$g''(x) = \frac{\exp(x) * (1 + \exp(x)) - \exp(x) * \exp(x)}{(1 + \exp(x))^2} = \frac{\exp(x)}{(1 + \exp(x))^2} \geq 0$$

Следовательно $g(x)$ является выпуклой.

Ф-ия $f(x) = -y \langle \mathbf{w}, \mathbf{x} \rangle$ является линейной.

Тогда ф-ия потерь будет выпуклой как композиция выпуклой ф-ии $f(x)$ и линейной ф-ии $g(x)$.

Следовательно, задача логистической регрессии является выпуклой задачей машинного обучения.

3. Докажите, что функция является выпуклой тогда и только тогда, когда её надграфик является выпуклым множеством

Пусть $f(x)$ выпуклая ф-ия. $f: \mathbb{R}^n \rightarrow \mathbb{R}$

Покажем, что её надграфик выпуклый.

Пусть (x_1, y_1) и (x_2, y_2) точки надграфика $f(x)$. Рассмотрим линию вида $(x', y') = a(x_1, y_1) + (1-a)(x_2, y_2)$, где $a \in [0, 1]$

Тогда

$y' = ay_1 + (1-a)y_2 \geq$ [т.к. это точки надграфика, то $y_1 \geq f(x_1)$]

$\geq af(x_1) + (1-a)f(x_2) \geq f(ax_1 + (1-a)x_2) = f(x')$

Т.е. y' лежит выше чем $f(x')$, значит, он содержится в надграфике (все точки линии лежат в надграфике). Следовательно надграфик это выпуклое множество.

Теперь в другую сторону.

Пусть надграфик выпуклый и x_1, x_2 из R^n , $a \in [0, 1]$

Тогда $(x_1, f(x_1))$ и $(x_2, f(x_2))$ находятся в надграфике $f(x)$, а так как он выпуклый то тогда $a(x_1, f(x_1)) + (1-a)(x_2, f(x_2))$ тоже лежит в надграфике. Тогда $f(ax_1 + (1-a)x_2) \leq af(x_1) + (1-a)f(x_2)$. Значит, $f(x)$ выпуклая.

- Докажите, что в случае разделимой выборки задача оптимизации функционала логистической регрессии не имеет решения, если не использовать регуляризацию

Функционал функции логистической регрессии имеет вид

$$\sum_{i=1}^{\ell} \ln(1 + \exp(-y_i \langle w, x_i \rangle)) \rightarrow \min_w.$$

Так как выборка разделима, то тогда существует вектор w такой что:

$\langle w, x_i \rangle > 0$ если $y_i = 1$ и $\langle w, x_i \rangle < 0$ если $y_i = -1$. Это и будет вектор разделяющий выборку (он может быть и не один). Тогда можно зафиксировать такой вектор и увеличивать его норму (умножая каждый элемент вектора на константу) и новый вектор также будет разделяющим, но при этом значение функционала можно будет постоянно уменьшать (оно будет как бы стремиться к нулю, но нуля мы достичь не сможем).

Т.е. в случае разделимой выборки у нас будет бесконечно много вариантов выбора вектора w . И решения задачи не будет существовать

А если будет регуляризация то выберем вектор с меньшей нормой.