

# **Численные методы проверки обоснованности обобщенных линейных моделей зависимостей**

**Левдик Павел Владимирович**  
Московский физико-технический институт

**Моттль Вадим Вячеславович**  
Вычислительный центр РАН

**Красоткина Ольга Вячеславовна**  
Московский государственный университет

**Татарчук Александр Игоревич**  
Вычислительный центр РАН

## Итак, мы рассматриваем вероятностную модель произвольной зависимости:

Априорная плотность распределения направляющего вектора и сдвига

$$\Psi(\mathbf{a}, b | \mu) \propto \exp\{-V(\mathbf{a} | \mu)\}$$

Условная плотность распределения характеристик случайного объекта

$$\varphi(\mathbf{x}, y | \mathbf{a}, b, c) \propto \exp\{-cq(y, z(\mathbf{x}, \mathbf{a}, b))\}$$

## Итак, мы рассматриваем вероятностную модель произвольной зависимости:

Регуляризирующая функция  
выбрана наблюдателем

$$V(\mathbf{a}|\mu) : \mathbb{R}^n \xrightarrow{\mu} \mathbb{R}^+$$

$$\Psi(\mathbf{a}, b|\mu) \propto \exp\{-V(\mathbf{a}|\mu)\}$$

Функция связи  
выбрана наблюдателем

$$q(y, z)$$

$$\varphi(\mathbf{x}, y/\mathbf{a}, b, c) \propto \exp\{-cq(y, z(\mathbf{x}, \mathbf{a}, b))\}$$

## Итак, мы рассматриваем вероятностную модель произвольной зависимости:

Регуляризирующая функция  
выбрана наблюдателем

$$V(\mathbf{a}|\mu) : \mathbb{R}^n \xrightarrow{\mu} \mathbb{R}^+$$

$$\Psi(\mathbf{a}, b|\mu) \propto \exp\{-V(\mathbf{a}|\mu)\}$$

Функция связи  
выбрана наблюдателем

$$q(y, z)$$

$$\varphi(\mathbf{x}, y/\mathbf{a}, b, c) \propto \exp\{-cq(y, z(\mathbf{x}, \mathbf{a}, b))\}$$

**Все остальное стандартно:**

Совместная условная плотность распределения  
случайной обучающей совокупности

$$\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) = \prod_{j=1}^N \varphi(\mathbf{x}_j, y_j / \mathbf{a}, b, c)$$

Структурные параметры обобщенной линейной  
модели зависимости

– параметр  $\mu$  регуляр. функции,  
– параметр  $c$  модели наблюдения.

## Итак, мы рассматриваем вероятностную модель произвольной зависимости:

Регуляризирующая функция  
выбрана наблюдателем

$$V(\mathbf{a}|\mu) : \mathbb{R}^n \xrightarrow{\mu} \mathbb{R}^+$$

$$\Psi(\mathbf{a}, b|\mu) \propto \exp\{-V(\mathbf{a}|\mu)\}$$

Функция связи  
выбрана наблюдателем

$$q(y, z)$$

$$\varphi(\mathbf{x}, y/\mathbf{a}, b, c) \propto \exp\{-cq(y, z(\mathbf{x}, \mathbf{a}, b))\}$$

**Все остальное стандартно:**

Совместная условная плотность распределения  
случайной обучающей совокупности

$$\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) = \prod_{j=1}^N \varphi(\mathbf{x}_j, y_j / \mathbf{a}, b, c)$$

Структурные параметры обобщенной линейной  
модели зависимости

– параметр  $\mu$  регуляр. функции,  
– параметр  $c$  модели наблюдения.

Маргинальная плотность распределения обучающей совокупности – непрерывная смесь условных распределений:

$$F(\mathbf{X}, \mathbf{Y} / \mu, c) = \int_{\mathbb{R}} \int_{\mathbb{R}^n} \Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b|\mu) da db$$

Функция правдоподобия для  $(\mu, c)$   
(Marginal Likelihood, Evidence Function)

Оценки максимального правдоподобия для структурных параметров:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(\mathbf{X}, \mathbf{Y} / \mu, c)$$

## Итак, мы рассматриваем вероятностную модель произвольной зависимости:

Регуляризирующая функция  
выбрана наблюдателем

$$V(\mathbf{a}|\mu) : \mathbb{R}^n \xrightarrow{\mu} \mathbb{R}^+$$

$$\Psi(\mathbf{a}, b|\mu) \propto \exp\{-V(\mathbf{a}|\mu)\}$$

Функция связи  
выбрана наблюдателем

$$q(y, z)$$

$$\varphi(\mathbf{x}, y/\mathbf{a}, b, c) \propto \exp\{-cq(y, z(\mathbf{x}, \mathbf{a}, b))\}$$

**Все остальное стандартно:**

Совместная условная плотность распределения  
случайной обучающей совокупности

$$\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) = \prod_{j=1}^N \varphi(\mathbf{x}_j, y_j / \mathbf{a}, b, c)$$

Структурные параметры обобщенной линейной  
модели зависимости

– параметр  $\mu$  регуляр. функции,  
– параметр  $c$  модели наблюдения.

Маргинальная плотность распределения обучающей совокупности – непрерывная смесь условных распределений:

$$F(\mathbf{X}, \mathbf{Y} / \mu, c) = \int_{\mathbb{R}} \int_{\mathbb{R}^n} \Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b|\mu) da db$$

Функция правдоподобия для  $(\mu, c)$   
(Marginal Likelihood, Evidence Function)

Оценки максимального правдоподобия для структурных параметров:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(\mathbf{X}, \mathbf{Y} / \mu, c)$$

Однако поиск точки максимума этой маргинальной плотности представляет собой крайне трудную вычислительную задачу.

## Альтернативная запись функции правдоподобия

Два эквивалентных выражения для совместная плотность распределения наблюдений, направляющего вектора и сдвига:

$$\underbrace{H(X, Y, a, b / \mu, c)}_{\text{совместное распределение}} = \underbrace{\Phi(X, Y / a, b, c)}_{\text{условное распределение}} \underbrace{\Psi(a, b / \mu)}_{\text{априорное распределение – функция правдоподобия}} = \underbrace{F(X, Y / \mu, c)}_{\text{маргинальное распределение}} \underbrace{P(a, b / X, Y, \mu, c)}_{\text{апостериорное распределение}}$$

## Альтернативная запись функции правдоподобия

Два эквивалентных выражения для совместная плотность распределения наблюдений, направляющего вектора и сдвига:

$$\underbrace{H(X, Y, a, b / \mu, c)}_{\text{совместное распределение}} = \underbrace{\Phi(X, Y / a, b, c)}_{\text{условное распределение}} \underbrace{\Psi(a, b | \mu)}_{\substack{\text{априорное} \\ \text{распределение} - \\ \text{функция} \\ \text{правдоподобия}}} = \underbrace{F(X, Y / \mu, c)}_{\text{маргинальное распределение}} \underbrace{P(a, b / X, Y, \mu, c)}_{\text{апостериорное распределение}}$$

Отсюда прямое выражение для функции правдоподобия:

При любых значениях  $(a, b)$

$$F(X, Y / \mu, c) = \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{P(a, b / X, Y, \mu, c)}$$



## Альтернативная запись функции правдоподобия

Два эквивалентных выражения для совместная плотность распределения наблюдений, направляющего вектора и сдвига:

$$\underbrace{H(X, Y, a, b / \mu, c)}_{\text{совместное распределение}} = \underbrace{\Phi(X, Y / a, b, c)}_{\text{условное распределение}} \underbrace{\Psi(a, b | \mu)}_{\substack{\text{априорное} \\ \text{распределение} - \\ \text{функция} \\ \text{правдоподобия}}} = \underbrace{F(X, Y / \mu, c)}_{\text{маргинальное распределение}} \underbrace{P(a, b / X, Y, \mu, c)}_{\text{апостериорное распределение}}$$

Отсюда прямое выражение для функции правдоподобия:

При любых значениях  $(a, b)$

$$F(X, Y / \mu, c) = \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{P(a, b / X, Y, \mu, c)}$$

Логарифмическая запись: Наиболее правдоподобные значения структурных параметров модели:

$$\ln F(X, Y / \mu, c) = \ln \Phi(X, Y / a, b, c) + \ln \Psi(a, b | \mu) - \ln P(a, b / X, Y, \mu, c) \rightarrow \max(\mu, c)$$

## Альтернативная запись функции правдоподобия

Два эквивалентных выражения для совместная плотность распределения наблюдений, направляющего вектора и сдвига:

$$\underbrace{H(X, Y, a, b / \mu, c)}_{\text{совместное распределение}} = \underbrace{\Phi(X, Y / a, b, c)}_{\text{условное распределение}} \underbrace{\Psi(a, b | \mu)}_{\substack{\text{априорное} \\ \text{распределение} - \\ \text{функция} \\ \text{правдоподобия}}} = \underbrace{F(X, Y / \mu, c)}_{\text{маргинальное распределение}} \underbrace{P(a, b / X, Y, \mu, c)}_{\text{апостериорное распределение}}$$

Отсюда прямое выражение для функции правдоподобия:

$$\text{При любых значениях } (a, b) \quad F(X, Y / \mu, c) = \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{P(a, b / X, Y, \mu, c)}$$

Логарифмическая запись: Наиболее правдоподобные значения структурных параметров модели:

$$\ln F(X, Y / \mu, c) = \ln \Phi(X, Y / a, b, c) + \ln \Psi(a, b | \mu) - \ln P(a, b / X, Y, \mu, c) \rightarrow \max(\mu, c)$$

Это представление логарифмической функции правдоподобия лежит в основе известного EM-алгоритма ее максимизации (Expectation- maximization).

## Нам нужно реализовать EM-алгоритм максимизации функции правдоподобия

Если бы мы могли вычислить апостериорную плотность для некоторых значений  $(\hat{\mu}_k, \hat{c}_k)$ , рассматриваемых как очередное приближение к точке максимума

$$\ln F(\mathbf{X}, \mathbf{Y} / \mu, c) = \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / a, b, c)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db +$$

$$\int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db - \underbrace{\int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \mu, c) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db}_{\text{negative term}},$$

то выполнялось бы  
фундаментальное неравенство:

$$\leq \int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db$$

## Нам нужно реализовать EM-алгоритм максимизации функции правдоподобия

Если бы мы могли вычислить апостериорную плотность для некоторых значений  $(\hat{\mu}_k, \hat{c}_k)$ , рассматриваемых как очередное приближение к точке максимума

$$\ln F(\mathbf{X}, \mathbf{Y} / \mu, c) = \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / a, b, c)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db +$$

$$\int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db - \underbrace{\int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \mu, c) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db}_{\text{expectation}},$$

то выполнялось бы фундаментальное неравенство:

$$\leq \int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db$$

Отсюда следует правило пересчета  $(\hat{\mu}_k, \hat{c}_k) \rightarrow (\hat{\mu}_{k+1}, \hat{c}_{k+1})$ , гарантирующее неубывание функции правдоподобия (EM- алгоритм):

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db,$$

$$\hat{c}_{k+1} = \underbrace{\arg \max_c}_{\text{maximization}} \underbrace{\int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / a, b, c)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db}_{\text{expectation}}.$$

## Нам нужно реализовать EM-алгоритм максимизации функции правдоподобия

Если бы мы могли вычислить апостериорную плотность для некоторых значений  $(\hat{\mu}_k, \hat{c}_k)$ , рассматриваемых как очередное приближение к точке максимума

$$\ln F(\mathbf{X}, \mathbf{Y} / \mu, c) = \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / a, b, c)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db + \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db - \underbrace{\int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \mu, c) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db}_{\text{expectation}}$$

то выполнялось бы фундаментальное неравенство:  $\leq \int \int_{\mathbb{R} \mathbb{R}^n} \ln P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db$

Отсюда следует правило пересчета  $(\hat{\mu}_k, \hat{c}_k) \rightarrow (\hat{\mu}_{k+1}, \hat{c}_{k+1})$ , гарантирующее неубывание функции правдоподобия (EM-алгоритм):

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db,$$

$$\hat{c}_{k+1} = \underbrace{\arg \max_c}_{\text{maximization}} \underbrace{\int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / a, b, c)] P(a, b / \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) da db}_{\text{expectation}}.$$

**Теорема:**  $\ln F(\mathbf{X}, \mathbf{Y} / \hat{\mu}_{k+1}, \hat{c}_{k+1}) \geq \ln F(\mathbf{X}, \mathbf{Y} / \hat{\mu}_k, \hat{c}_k)$ ,

пока  $\frac{\partial}{\partial \mu} \ln F(\mathbf{X}, \mathbf{Y} / \hat{\mu}_k, \hat{c}_k) \neq 0$  либо  $\frac{\partial}{\partial c} \ln F(\mathbf{X}, \mathbf{Y} / \hat{\mu}_k, \hat{c}_k) \neq 0$ .

## Численно реализовать исходный EM-алгоритм невозможно

Еще раз EM-алгоритм:

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db,$$

$$\hat{c}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db.$$

## Численно реализовать исходный EM-алгоритм невозможно

Еще раз EM-алгоритм:

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db,$$

$$\hat{c}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db.$$

Апостериорное распределение направляющего вектора и сдвига:

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \propto \Psi(\mathbf{a}, b | \mu) \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)$$

## Численно реализовать исходный EM-алгоритм невозможно

Еще раз EM-алгоритм:

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db,$$

$$\hat{c}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db.$$

Апостериорное распределение направляющего вектора и сдвига:

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \propto \Psi(\mathbf{a}, b | \mu) \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)$$

$$\Psi(\mathbf{a}, b | \mu) \propto \exp\{-V(\mathbf{a} | \mu)\}, \quad \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c) \propto \exp\left\{-c \sum_{j=1}^N q(y_j, z(\mathbf{x}_j, \mathbf{a}, b))\right\}$$

Для многих регуляризирующих функций  $V(\mathbf{a} | \mu)$  и функций связи  $q(y, z)$  интегрирование по такой апостериорной плотности затруднительно.



## Численно реализовать исходный EM-алгоритм невозможно

Еще раз EM-алгоритм:

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db,$$

$$\hat{c}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db.$$

Апостериорное распределение направляющего вектора и сдвига:

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \propto \Psi(\mathbf{a}, b | \mu) \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)$$

$$\Psi(\mathbf{a}, b | \mu) \propto \exp\{-V(\mathbf{a} | \mu)\}, \quad \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c) \propto \exp\left\{-c \sum_{j=1}^N q(y_j, z(\mathbf{x}_j, \mathbf{a}, b))\right\}$$

Для многих регуляризирующих функций  $V(\mathbf{a} | \mu)$  и функций связи  $q(y, z)$  интегрирование по такой апостериорной плотности затруднительно.

## Идея: Нормальная аппроксимация апостериорной плотности

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \cong \underbrace{\mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c})}_{\text{маргинальное апостериорное распределение}} \underbrace{\mathcal{N}(b | \hat{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \hat{\sigma}_N^2)}_{\text{условное апостериорное распределение}}.$$

# Численно реализовать исходный EM-алгоритм невозможно

Еще раз EM-алгоритм:

$$\hat{\mu}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db,$$

$$\hat{c}_{k+1} = \arg \max_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)] P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \hat{\mu}_k, \hat{c}_k) d\mathbf{a} db.$$

Апостериорное распределение направляющего вектора и сдвига:

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \propto \Psi(\mathbf{a}, b | \mu) \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c)$$

$$\Psi(\mathbf{a}, b | \mu) \propto \exp\{-V(\mathbf{a} | \mu)\}, \quad \Phi(\mathbf{X}, \mathbf{Y} | \mathbf{a}, b, c) \propto \exp\left\{-c \sum_{j=1}^N q(y_j, z(\mathbf{x}_j, \mathbf{a}, b))\right\}$$

Для многих регуляризующих функций  $V(\mathbf{a} | \mu)$  и функций связи  $q(y, z)$  интегрирование по такой апостериорной плотности затруднительно.

## Идея: Нормальная аппроксимация апостериорной плотности

$$P(\mathbf{a}, b | \mathbf{X}, \mathbf{Y}, \mu, c) \cong \underbrace{\mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c})}_{\text{маргинальное апостериорное распределение}} \underbrace{\mathcal{N}(b | \hat{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \hat{\sigma}_N^2)}_{\text{условное апостериорное распределение}}.$$

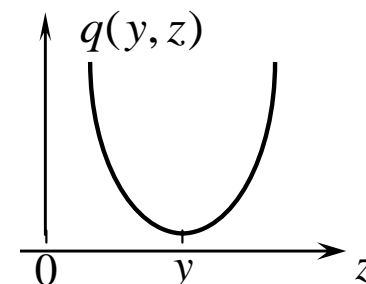
Ограничимся задачей восстановления числовой регрессии с квадратично-модульной регуляризацией

Регуляризация

Elastic Net  $\mathbf{a} \in \mathbb{R}^n$

$$V(\mathbf{a} | \mu) = \sum_{i=1}^n a_i^2 + \mu \sum_{i=1}^n |a_i|$$

Квадратичная функция связи



## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

**Теорема:** Пусть  $(\hat{\delta}_1, \dots, \hat{\delta}_N)$  есть решение задачи выпуклой оптимизации

$$W_{y, X, \mu, c}(\delta_1, \dots, \delta_N) = \sum_{i=1}^n \left\{ \min \left[ \frac{\mu}{2} + \sum_{j=1}^N \delta_j x_{ij}, 0, \frac{\mu}{2} - \sum_{j=1}^N \delta_j x_{ij} \right] \right\}^2 + 2c \sum_{j=1}^N (\delta_j - y_j)^2 \rightarrow \begin{cases} \min(\delta_1, \dots, \delta_N), \\ \sum_{j=1}^N \delta_j = 0. \end{cases}$$

## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

**Теорема:** Пусть  $(\hat{\delta}_1, \dots, \hat{\delta}_N)$  есть решение задачи выпуклой оптимизации

$$W_{y, X, \mu, c}(\delta_1, \dots, \delta_N) = \sum_{i=1}^n \left\{ \min \left[ \frac{\mu}{2} + \sum_{j=1}^N \delta_j x_{ij}, 0, \frac{\mu}{2} - \sum_{j=1}^N \delta_j x_{ij} \right] \right\}^2 + 2c \sum_{j=1}^N (\delta_j - y_j)^2 \rightarrow \begin{cases} \min(\delta_1, \dots, \delta_N), \\ \sum_{j=1}^N \delta_j = 0. \end{cases}$$

Тогда оценки параметров регрессии определяются выражениями:

$$\hat{a}_{i, y, X, \mu, c} = \frac{1}{\beta} \begin{cases} \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} + \mu \right) < 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} < -\frac{\mu}{2}, \\ 0, & -\frac{\mu}{2} \leq \sum_{j=1}^N \hat{\delta}_j x_{ij} \leq \frac{\mu}{2}, \\ \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} - \mu \right) > 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} > \frac{\mu}{2}, \end{cases} \quad \begin{matrix} i = 1, \dots, n, \\ \hat{\mathbf{a}}_{y, X, \mu, c} = (\hat{a}_{1, y, X, \mu, c} \cdots \hat{a}_{n, y, X, \mu, c})^T \in \mathbb{R}^n, \end{matrix}$$

$$\hat{b}_{y, X, \mu, c} = \frac{1}{N} \sum_{j=1}^N \left( y_j - \sum_{i=1}^n x_{ij} \hat{a}_{i, y, X, \mu, c} \right).$$

## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

**Теорема:** Пусть  $(\hat{\delta}_1, \dots, \hat{\delta}_N)$  есть решение задачи выпуклой оптимизации

$$W_{y, X, \mu, c}(\delta_1, \dots, \delta_N) = \sum_{i=1}^n \left\{ \min \left[ \frac{\mu}{2} + \sum_{j=1}^N \delta_j x_{ij}, 0, \frac{\mu}{2} - \sum_{j=1}^N \delta_j x_{ij} \right] \right\}^2 + 2c \sum_{j=1}^N (\delta_j - y_j)^2 \rightarrow \begin{cases} \min(\delta_1, \dots, \delta_N), \\ \sum_{j=1}^N \delta_j = 0. \end{cases}$$

Тогда оценки параметров регрессии определяются выражениями:

$$\hat{a}_{i, y, X, \mu, c} = \frac{1}{\beta} \begin{cases} \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} + \mu \right) < 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} < -\frac{\mu}{2}, \\ 0, & -\frac{\mu}{2} \leq \sum_{j=1}^N \hat{\delta}_j x_{ij} \leq \frac{\mu}{2}, \\ \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} - \mu \right) > 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} > \frac{\mu}{2}, \end{cases} \quad \begin{matrix} i = 1, \dots, n, \\ \hat{\mathbf{a}}_{y, X, \mu, c} = (\hat{a}_{1, y, X, \mu, c} \cdots \hat{a}_{n, y, X, \mu, c})^T \in \mathbb{R}^n, \end{matrix}$$

$$\hat{b}_{y, X, \mu, c} = \frac{1}{N} \sum_{j=1}^N \left( y_j - \sum_{i=1}^n x_{ij} \hat{a}_{i, y, X, \mu, c} \right).$$

Очевидный эффект автоматического отбора признаков.

## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

**Теорема:** Пусть  $(\hat{\delta}_1, \dots, \hat{\delta}_N)$  есть решение задачи выпуклой оптимизации

$$W_{y, X, \mu, c}(\delta_1, \dots, \delta_N) = \sum_{i=1}^n \left\{ \min \left[ \frac{\mu}{2} + \sum_{j=1}^N \delta_j x_{ij}, 0, \frac{\mu}{2} - \sum_{j=1}^N \delta_j x_{ij} \right] \right\}^2 + 2c \sum_{j=1}^N (\delta_j - y_j)^2 \rightarrow \begin{cases} \min(\delta_1, \dots, \delta_N), \\ \sum_{j=1}^N \delta_j = 0. \end{cases}$$

Тогда оценки параметров регрессии определяются выражениями:

$$\hat{a}_{i, y, X, \mu, c} = \frac{1}{\beta} \begin{cases} \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} + \mu \right) < 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} < -\frac{\mu}{2}, \\ 0, & -\frac{\mu}{2} \leq \sum_{j=1}^N \hat{\delta}_j x_{ij} \leq \frac{\mu}{2}, \\ \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} - \mu \right) > 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} > \frac{\mu}{2}, \end{cases} \quad i = 1, \dots, n,$$

$$\hat{\mathbf{a}}_{y, X, \mu, c} = (\hat{a}_{1, y, X, \mu, c} \cdots \hat{a}_{n, y, X, \mu, c})^T \in \mathbb{R}^n,$$

$$\hat{b}_{y, X, \mu, c} = \frac{1}{N} \sum_{j=1}^N \left( y_j - \sum_{i=1}^n x_{ij} \hat{a}_{i, y, X, \mu, c} \right).$$

Очевидный эффект автоматического отбора признаков.

Три подмножества компонент вектора коэффициентов регрессии:

$$\hat{\mathbb{I}}_{y, X, \mu, c}^+ = \{i: \hat{a}_{i, y, X, \mu, c} > 0\}, \quad \hat{\mathbb{I}}_{y, X, \mu, c}^- = \{i: \hat{a}_{i, y, X, \mu, c} < 0\}$$

$$\hat{\mathbb{I}}_{y, X, \mu, c}^0 = \{i: \hat{a}_{i, y, X, \mu, c} = 0\}$$

## Оценивание параметров регрессии: Минимум регуляризованного эмпирического риска

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$(\hat{\mathbf{a}}_{y, X, \mu, c}, \hat{b}_{y, X, \mu, c}) = \arg \min_{\mathbf{a}, b} \left\{ \sum_{i=1}^n (a_i^2 + \mu |a_i|) + \frac{1}{2c} \sum_{j=1}^N \left[ y_j - \left( \sum_{i=1}^n x_{ij} a_i + b \right) \right]^2 \right\}$$

**Теорема:** Пусть  $(\hat{\delta}_1, \dots, \hat{\delta}_N)$  есть решение задачи выпуклой оптимизации

$$W_{y, X, \mu, c}(\delta_1, \dots, \delta_N) = \sum_{i=1}^n \left\{ \min \left[ \frac{\mu}{2} + \sum_{j=1}^N \delta_j x_{ij}, 0, \frac{\mu}{2} - \sum_{j=1}^N \delta_j x_{ij} \right] \right\}^2 + 2c \sum_{j=1}^N (\delta_j - y_j)^2 \rightarrow \begin{cases} \min(\delta_1, \dots, \delta_N), \\ \sum_{j=1}^N \delta_j = 0. \end{cases}$$

Тогда оценки параметров регрессии определяются выражениями:

$$\hat{a}_{i, y, X, \mu, c} = \frac{1}{\beta} \begin{cases} \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} + \mu \right) < 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} < -\frac{\mu}{2}, \\ 0, & -\frac{\mu}{2} \leq \sum_{j=1}^N \hat{\delta}_j x_{ij} \leq \frac{\mu}{2}, \\ \left( \sum_{j=1}^N \hat{\delta}_j x_{ij} - \mu \right) > 0, & \sum_{j=1}^N \hat{\delta}_j x_{ij} > \frac{\mu}{2}, \end{cases} \quad i = 1, \dots, n,$$

$$\hat{\mathbf{a}}_{y, X, \mu, c} = (\hat{a}_{1, y, X, \mu, c} \cdots \hat{a}_{n, y, X, \mu, c})^T \in \mathbb{R}^n,$$

$$\hat{b}_{y, X, \mu, c} = \frac{1}{N} \sum_{j=1}^N \left( y_j - \sum_{i=1}^n x_{ij} \hat{a}_{i, y, X, \mu, c} \right).$$

Очевидный эффект автоматического отбора признаков.

Чем больше параметр селективности  $\mu \geq 0$ , тем меньше ненулевых коэффициентов регрессии, т.е. меньше уровень сложности модели.



## Нормальная аппроксимация апостериорной плотности

Исходная «ненормальная» апостериорная плотность случайных параметров регрессии:

$$P(\mathbf{a}, b \mid \mathbf{y}, \mathbf{X}, \mu, c) \propto \Psi(\mathbf{a} \mid \mu) \Phi(\mathbf{y} \mid \mathbf{X}, \mathbf{a}, b, c) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^0} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} |a_i| + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

## Нормальная аппроксимация апостериорной плотности

Исходная «ненормальная» апостериорная плотность случайных параметров регрессии:

$$P(\mathbf{a}, \mathbf{b} | \mathbf{y}, \mathbf{X}, \mu, c) \propto \Psi(\mathbf{a} | \mu) \Phi(\mathbf{y} | \mathbf{X}, \mathbf{a}, \mathbf{b}, c) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^0} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} |a_i| + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

Нормальная аппроксимация:

$$\hat{P}(\mathbf{a}, \mathbf{b} | \mathbf{y}, \mathbf{X}, \mu, c) = \mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{\mathbf{y}, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{\mathbf{y}, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} a_i - \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} a_i + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

## Нормальная аппроксимация апостериорной плотности

Исходная «ненормальная» апостериорная плотность случайных параметров регрессии:

$$P(\mathbf{a}, \mathbf{b} | \mathbf{y}, \mathbf{X}, \mu, c) \propto \Psi(\mathbf{a} | \mu) \Phi(\mathbf{y} | \mathbf{X}, \mathbf{a}, \mathbf{b}, c) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^0} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} |a_i| + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

Нормальная аппроксимация:

$$\hat{P}(\mathbf{a}, \mathbf{b} | \mathbf{y}, \mathbf{X}, \mu, c) = \mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{\mathbf{y}, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{\mathbf{y}, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} \quad \quad \quad + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} a_i - \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} a_i + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

## Нормальная аппроксимация апостериорной плотности

Исходная «ненормальная» апостериорная плотность случайных параметров регрессии:

$$P(\mathbf{a}, b | \mathbf{y}, \mathbf{X}, \mu, c) \propto \Psi(\mathbf{a} | \mu) \Phi(\mathbf{y} | \mathbf{X}, \mathbf{a}, b, c) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^0} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} |a_i| + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} |a_i| + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

Нормальная аппроксимация:

$$\hat{P}(\mathbf{a}, b | \mathbf{y}, \mathbf{X}, \mu, c) = \mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{\mathbf{y}, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{\mathbf{y}, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2) \propto$$

$$\exp \left\{ - \left\{ \mathbf{a}^T \mathbf{a} + \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^+} a_i - \mu \sum_{i \in \hat{\mathbb{I}}_{\mathbf{y}, \mathbf{X}, \mu, c}^-} a_i + \frac{1}{2c} [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)]^T [\mathbf{y} - (\mathbf{X}^T \mathbf{a} + \mathbf{1}b)] \right\} \right\}.$$

Можно доказать, что здесь:

апостериорная ковариационная матрица  $\mathbf{B}_{\mathbf{X}, c} = \left( \frac{1}{c} \mathbf{X} \mathbf{X}^T + 2\mathbf{I} \right)^{-1},$

апостериорное условное математическое ожидание  $\tilde{b}_{\mathbf{y}, \mathbf{X}, \mu, c}(\mathbf{a}) = \hat{b}_{\mathbf{y}, \mathbf{X}, \mu, c} - \frac{1}{2Nc} \mathbf{1}^T \mathbf{X}^T (\mathbf{a} - \hat{\mathbf{a}}_{\mathbf{y}, \mathbf{X}, \mu, c}),$

апостериорная условная дисперсия  $\tilde{\sigma}_N^2 = \frac{1}{2N}.$

## Приближенная функция правдоподобия

Исходный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(X, Y / \mu, c) = \operatorname{argmax} \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{P(a, b / X, Y, \mu, c)}.$$

## Приближенная функция правдоподобия

Исходный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(\mathbf{X}, \mathbf{Y} / \mu, c) = \operatorname{argmax} \frac{\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b | \mu)}{P(\mathbf{a}, b / \mathbf{X}, \mathbf{Y}, \mu, c)}.$$

Приближенный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} \hat{F}(\mathbf{X}, \mathbf{Y} / \mu, c) = \operatorname{argmax} \frac{\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b | \mu)}{\mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2)}.$$

## Приближенная функция правдоподобия

Исходный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(\mathbf{X}, \mathbf{Y} / \mu, c) = \operatorname{argmax} \frac{\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b | \mu)}{P(\mathbf{a}, b / \mathbf{X}, \mathbf{Y}, \mu, c)}.$$

Приближенный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} \hat{F}(\mathbf{X}, \mathbf{Y} / \mu, c) = \operatorname{argmax} \frac{\Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c) \Psi(\mathbf{a}, b | \mu)}{\mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2)}.$$

EM-алгоритм для приближенного критерия максимума правдоподобия:

$$\begin{aligned} \hat{\mu}_{k+1} &= \operatorname{argmax}_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(\mathbf{a}, b | \mu)] \mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2) d\mathbf{a} db, \\ \hat{c}_{k+1} &= \operatorname{argmax}_c \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(\mathbf{X}, \mathbf{Y} / \mathbf{a}, b, c)] \mathcal{N}(\mathbf{a} | \hat{\mathbf{a}}_{y, \mathbf{X}, \mu, c}, \mathbf{B}_{\mathbf{X}, c}) \mathcal{N}(b | \tilde{b}_{y, \mathbf{X}, \mu, c}(\mathbf{a}), \tilde{\sigma}_N^2) d\mathbf{a} db. \end{aligned}$$

## Приближенная функция правдоподобия

Исходный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} F(X, Y / \mu, c) = \operatorname{argmax} \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{P(a, b / X, Y, \mu, c)}.$$

Приближенный критерий максимума правдоподобия:

$$(\hat{\mu}, \hat{c}) = \operatorname{argmax} \hat{F}(X, Y / \mu, c) = \operatorname{argmax} \frac{\Phi(X, Y / a, b, c) \Psi(a, b | \mu)}{\mathcal{N}(a | \hat{a}_{y, X, \mu, c}, \mathbf{B}_{X, c}) \mathcal{N}(b | \tilde{b}_{y, X, \mu, c}(a), \tilde{\sigma}_N^2)}.$$

EM-алгоритм для приближенного критерия максимума правдоподобия:

$$\begin{aligned} \hat{\mu}_{k+1} &= \operatorname{arg max}_{\mu} \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Psi(a, b | \mu)] \mathcal{N}(a | \hat{a}_{y, X, \mu, c}, \mathbf{B}_{X, c}) \mathcal{N}(b | \tilde{b}_{y, X, \mu, c}(a), \tilde{\sigma}_N^2) da db, \\ \hat{c}_{k+1} &= \operatorname{arg max}_c \int \int_{\mathbb{R} \mathbb{R}^n} [\ln \Phi(X, Y / a, b, c)] \mathcal{N}(a | \hat{a}_{y, X, \mu, c}, \mathbf{B}_{X, c}) \mathcal{N}(b | \tilde{b}_{y, X, \mu, c}(a), \tilde{\sigma}_N^2) da db. \end{aligned}$$

**Численная реализация такого алгоритма не составляет проблемы.**



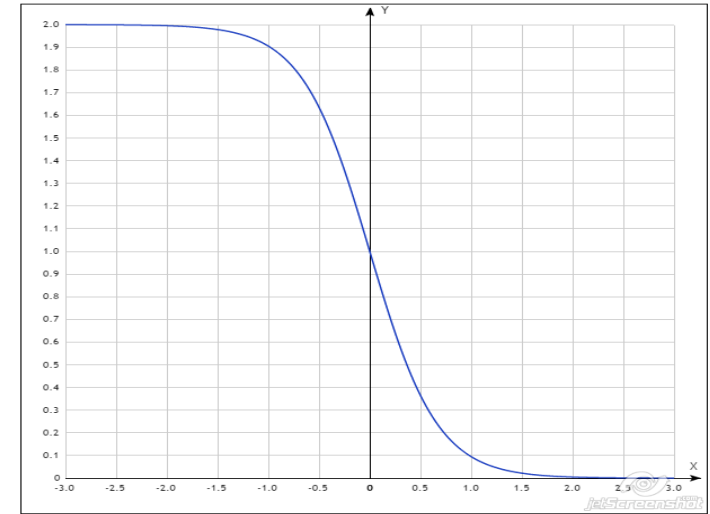
# Приближенная функция правдоподобия

## Шаг по параметру селективности

$$\hat{\mu}_{k+1} = \arg \min_{\mu \geq 0} \left[ \underbrace{-\frac{1}{\sqrt{\pi}} \exp\left(-\frac{\mu^2}{4}\right) \left[\operatorname{erfc}\left(\frac{\mu}{2}\right)\right]^{-1}}_{\text{выпуклая функция}} + h_k \mu \right]$$

Здесь:  $\operatorname{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^{\infty} \exp(-t^2) dt$  – функция ошибок,

$h_k(\hat{a}_{i,y,X}, \hat{\mu}_k, \hat{c}_k, \tilde{\sigma}_N^2)$  – константа, определяемая аппроксимирующим нормальным распределением.



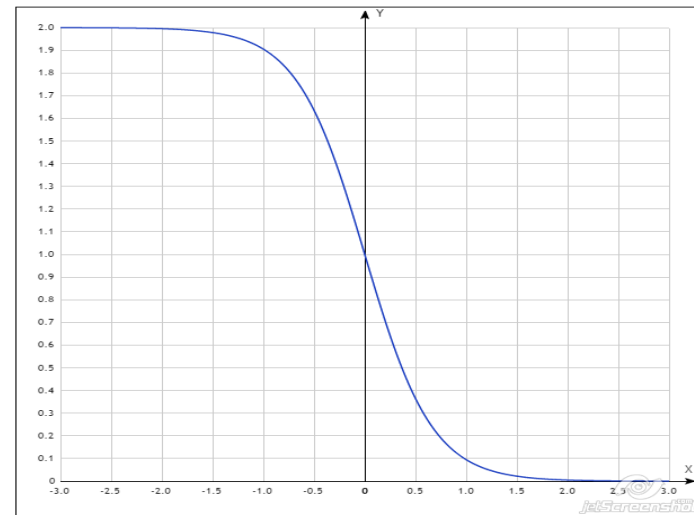
# Приближенная функция правдоподобия

## Шаг по параметру селективности

$$\hat{\mu}_{k+1} = \arg \min_{\mu \geq 0} \underbrace{\left[ -\frac{1}{\sqrt{\pi}} \exp\left(-\frac{\mu^2}{4}\right) \left[ \operatorname{erfc}\left(\frac{\mu}{2}\right) \right]^{-1} + h_k \mu \right]}_{\text{выпуклая функция}}$$

Здесь:  $\operatorname{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^{\infty} \exp(-t^2) dt$  – функция ошибок,

$h_k(\hat{a}_{i,y,X,\hat{\mu}_k,\hat{c}_k}, \tilde{\sigma}_N^2)$  – константа, определяемая аппроксимирующим нормальным распределением.



## Шаг по параметру наблюдений

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j=1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$\hat{c}_{k+1} = \frac{1}{N} \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right)^T \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right) + \frac{1}{4\hat{c}_k} \operatorname{Tr} \left[ \left( \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \mathbf{x}_j \mathbf{x}_l^T \right) (\mathbf{X} \mathbf{X}^T + 2\hat{c}_k \mathbf{I})^{-1} \right] + \tilde{\sigma}_N^2.$$

Здесь  $(\hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k}, \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k})$  – результат обучения на предыдущем шаге.

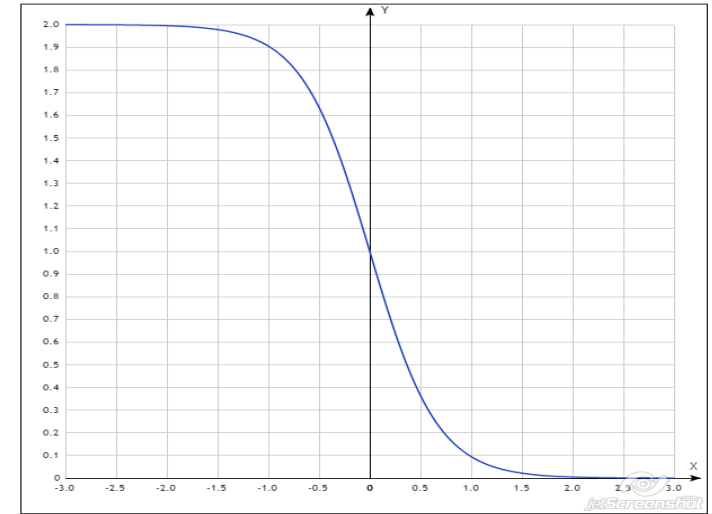
# Приближенная функция правдоподобия

## Шаг по параметру селективности

$$\hat{\mu}_{k+1} = \arg \min_{\mu \geq 0} \underbrace{\left[ -\frac{1}{\sqrt{\pi}} \exp\left(-\frac{\mu^2}{4}\right) \left[ \operatorname{erfc}\left(\frac{\mu}{2}\right) \right]^{-1} + h_k \mu \right]}_{\text{выпуклая функция}}$$

Здесь:  $\operatorname{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^{\infty} \exp(-t^2) dt$  – функция ошибок,

$h_k(\hat{a}_{i,y,X,\hat{\mu}_k,\hat{c}_k}, \tilde{\sigma}_N^2)$  – константа, определяемая аппроксимирующим нормальным распределением.



## Шаг по параметру наблюдений

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$\hat{c}_{k+1} = \frac{1}{N} \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right)^T \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right) + \frac{1}{4\hat{c}_k} \operatorname{Tr} \left[ \left( \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \mathbf{x}_j \mathbf{x}_l^T \right) (\mathbf{X} \mathbf{X}^T + 2\hat{c}_k \mathbf{I})^{-1} \right] + \tilde{\sigma}_N^2.$$

Здесь  $(\hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k}, \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k})$  – результат обучения на предыдущем шаге.

Алгоритм не требует перебора значений структурных параметров  $(\mu, c)$ , как в известном методе Regularization Path.

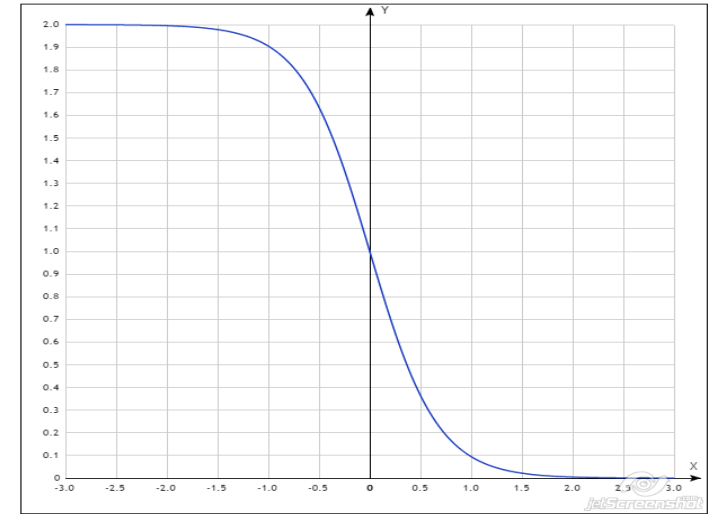
# Приближенная функция правдоподобия

## Шаг по параметру селективности

$$\hat{\mu}_{k+1} = \arg \min_{\mu \geq 0} \underbrace{\left[ -\frac{1}{\sqrt{\pi}} \exp\left(-\frac{\mu^2}{4}\right) \left[ \operatorname{erfc}\left(\frac{\mu}{2}\right) \right]^{-1} + h_k \mu \right]}_{\text{выпуклая функция}}$$

Здесь:  $\operatorname{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^{\infty} \exp(-t^2) dt$  – функция ошибок,

$h_k(\hat{a}_{i,y,X,\hat{\mu}_k,\hat{c}_k}, \tilde{\sigma}_N^2)$  – константа, определяемая аппроксимирующим нормальным распределением.



## Шаг по параметру наблюдений

Обучающая совокупность  $\{(\mathbf{x}_j, y_j), j = 1, \dots, N\}$ ,  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $y_j \in \mathbb{R}$

$$\hat{c}_{k+1} = \frac{1}{N} \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right)^T \left( \mathbf{y} - (\mathbf{X}^T \hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k} + \mathbf{I} \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k}) \right) + \frac{1}{4\hat{c}_k} \operatorname{Tr} \left[ \left( \frac{1}{N^2} \sum_{j=1}^N \sum_{l=1}^N \mathbf{x}_j \mathbf{x}_l^T \right) (\mathbf{X} \mathbf{X}^T + 2\hat{c}_k \mathbf{I})^{-1} \right] + \tilde{\sigma}_N^2.$$

Здесь  $(\hat{\mathbf{a}}_{y,X,\hat{\mu}_k,\hat{c}_k}, \hat{\mathbf{b}}_{y,X,\hat{\mu}_k,\hat{c}_k})$  – результат обучения на предыдущем шаге.

J. Friedman, T. Hastie, R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 2010, Volume 33, Issue 1.