

# Лекция 7. Недиагональная регуляризация обобщенных линейных моделей

Д. П. Ветров<sup>1</sup>    Д. А. Кропотов<sup>2</sup>

<sup>1</sup>МГУ, ВМиК, каф. ММП

<sup>2</sup>ВЦ РАН

Спецкурс «Байесовские методы машинного обучения»

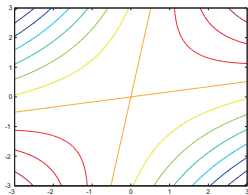
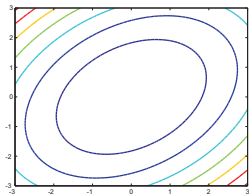
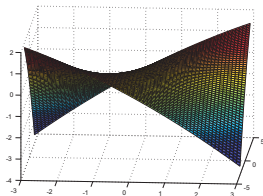
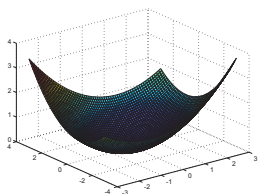
# План лекции

Ликбез

Метод релевантных собственных векторов

# Неотрицательно определенные матрицы I

$$A \in \mathbb{R}^{n \times n}, A = A^T \leftrightarrow f(\mathbf{x}) = \langle A\mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T A \mathbf{x}$$



Канонический вид квадратичной формы:

$$\exists P : P^T = P^{-1}, \det(P) \neq 0, \mathbf{y} = P\mathbf{x}, f(\mathbf{y}) = \sum_{j=1}^n \lambda_j y_j^2, A = \text{diag}\{\lambda_1, \dots, \lambda_n\}$$

# Неотрицательно определенные матрицы II

Матрица — неотрицательно определена, если соответствующая ей квадратичная форма всегда неотрицательна:

$$\forall \mathbf{x} \in \mathbb{R}^n \quad \mathbf{x}^T A \mathbf{x} \geq 0 \text{ и } = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}.$$

У неотрицательно определенной матрицы  $A$  все собственные значения  $\lambda_j \geq 0$

Количество ненулевых  $\lambda_j$  называется рангом матрицы  $A$ .

# Линейная регрессия I

Выборка  $(X, \mathbf{t}) = \{\mathbf{x}_n, t_n\}_{n=1}^N$ ,  $\mathbf{x}_n \in \mathbb{R}^d$ ,  $t_n \in \mathbb{R}$

Решающее правило:

$$y(\mathbf{x}) = \sum_{j=1}^K w_j \phi_j(\mathbf{x})$$

Обучение в детерминистской постановке:

$$\begin{aligned} \sum_{n=1}^N (t_n - y(\mathbf{x}_n))^2 + \lambda \|\mathbf{w}\|^2 &= \sum_{n=1}^N \left( t_n - \sum_{j=1}^K w_j \phi_j(\mathbf{x}_n) \right)^2 + \lambda \|\mathbf{w}\|^2 = \\ &= \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}} \end{aligned}$$

$$\mathbf{w}_{opt} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{t}$$

# Линейная регрессия II

Обучение в вероятностной постановке:

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}), \sigma^2), \quad p(\mathbf{t}|X, \mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}), \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \lambda I)$$

$$p(\mathbf{w}|X, \mathbf{t}) \rightarrow \max_{\mathbf{w}} \Leftrightarrow p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w}) \rightarrow \max_{\mathbf{w}} \Leftrightarrow$$

$$\log p(\mathbf{t}|X, \mathbf{w}) + \log p(\mathbf{w}) \rightarrow \max_{\mathbf{w}} \Leftrightarrow \frac{1}{2\sigma^2} \|\mathbf{t} - \Phi\mathbf{w}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}$$

$$\mathbf{w}_{MP} = \left( \frac{1}{\sigma^2} \Phi^T \Phi + \lambda I \right)^{-1} \frac{1}{\sigma^2} \Phi^T \mathbf{t}$$

# Логистическая регрессия

Выборка  $(X, \mathbf{t}) = \{\mathbf{x}_n, t_n\}_{n=1}^N$ ,  $\mathbf{x}_n \in \mathbb{R}^d$ ,  $t_n \in \{-1, 1\}$

Решающее правило:

$$y(\mathbf{x}) = \text{sign}(f(\mathbf{x})) = \text{sign}\left(\sum_{j=1}^K w_j \phi_j(\mathbf{x})\right)$$

Обучение в детерминистской постановке:

$$\sum_{n=1}^N \log(1 + \exp(-t_n f(\mathbf{x}_n))) + \lambda \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}$$

Обучение в вероятностной постановке:

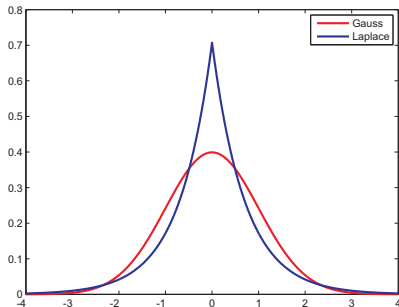
$$p(t|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-t f(\mathbf{x}))}, \quad p(\mathbf{t}|X, \mathbf{w}) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}), \quad p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, \lambda I)$$

$$p(\mathbf{w}|X, \mathbf{t}) \rightarrow \max_{\mathbf{w}} \Leftrightarrow \sum_{n=1}^N \log(1 + \exp(-t_n f(\mathbf{x}_n))) + \frac{\lambda}{2} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}}$$

# Распределение Лапласа I

Плотность распределения:

$$\mathcal{L}(x|\lambda) = \frac{\lambda}{4} \exp\left(-\frac{\lambda}{2}|x|\right)$$



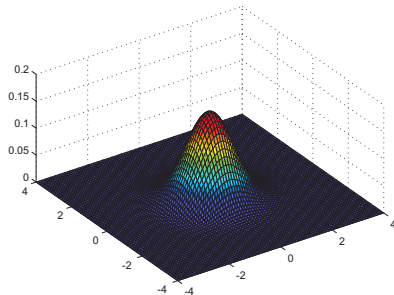
Распределение Лапласа имеет сингулярность в нуле и более тяжелые хвосты, чем нормальное распределение



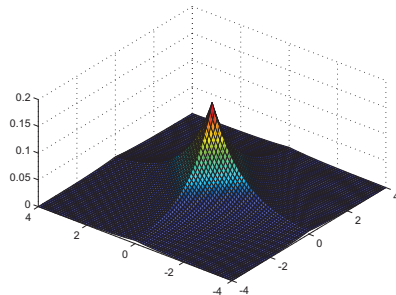
# Распределение Лапласа II

Многомерное распределение Лапласа:

$$\mathcal{L}(\mathbf{x}|\lambda) = \left(\frac{\lambda}{4}\right)^d \exp\left(-\frac{\lambda}{2} \sum_{j=1}^d |x_j|\right)$$



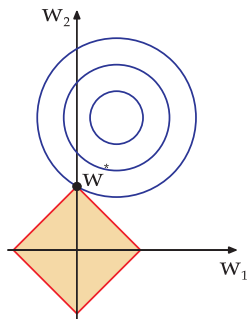
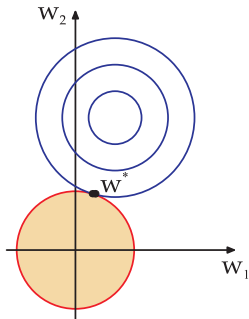
Нормальное распределение  
 $\mathcal{N}(\mathbf{x}|\lambda) \propto \exp\left(-\frac{\lambda}{2} \|\mathbf{x}\|_{\ell_2}^2\right)$



Распределение Лапласа  
 $\mathcal{L}(\mathbf{x}|\lambda) \propto \exp\left(-\frac{\lambda}{2} \|\mathbf{x}\|_{\ell_1}\right)$

# Лапласовская регуляризация

$$\log p(\mathbf{t}|\mathbf{w}, X) - \frac{\alpha}{2} \sum_{j=1}^K |w_j| \rightarrow \max_{\mathbf{w}} \leftrightarrow \log p(\mathbf{t}|\mathbf{w}, X) \rightarrow \max_{\mathbf{w}} \\ -\frac{1}{2} \sum_{j=1}^K |w_j| \leq \gamma$$



# Метод релевантных векторов

Регуляризатор:

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \sqrt{\frac{\det(A)}{(2\pi)^K}} \exp\left(-\frac{1}{2}\mathbf{w}^T A \mathbf{w}\right),$$

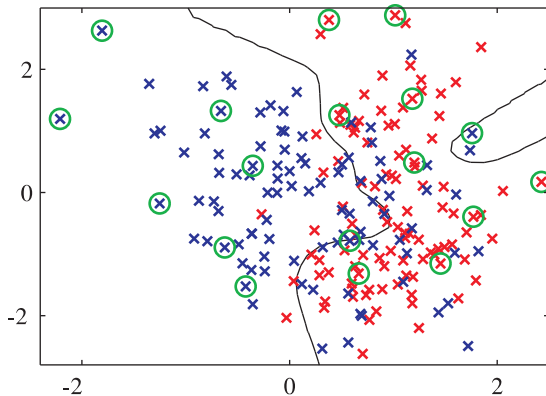
где  $A = \text{diag}(\alpha_1, \dots, \alpha_K)$ ,  $\alpha_j \geq 0$ .

Для определения значений коэффициентов регуляризации используется принцип наибольшей обоснованности

$$\boldsymbol{\alpha}_{ME} = \arg \max_{\boldsymbol{\alpha}} p(\mathbf{t}|X, \boldsymbol{\alpha}) = \arg \max_{\boldsymbol{\alpha}} \int p(\mathbf{t}|X, \mathbf{w}) p(\mathbf{w}|\boldsymbol{\alpha}) d\mathbf{w}.$$

# Разреженность RVM

В процессе оптимизации обоснованности большинство  $\alpha_j$  уходят в  $+\infty$ , исключая нерелевантные базисные функции из решающего правила.



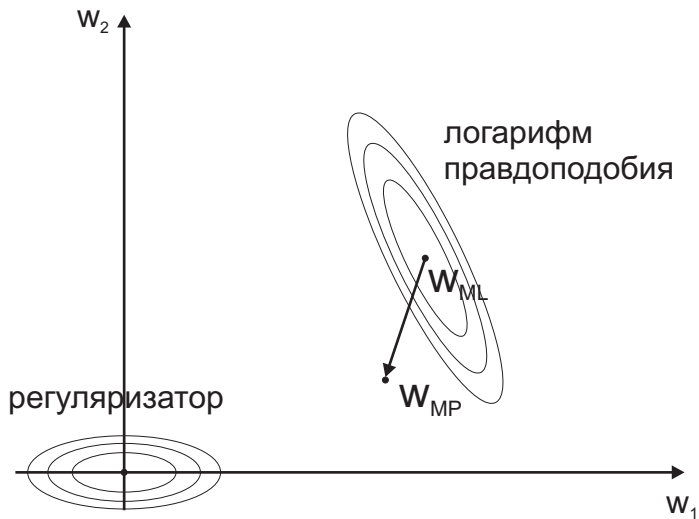
# Недостатки RVM

- При обучении классификатора RVM требуется порядка 20–50 итераций для настройки  $\alpha$ , на каждой из которых приходится обучать метод регуляризованной линейной/логистической регрессии
- RVM не может быть напрямую применен для лапласовского априорного распределения на веса  $w$ . В то же время известно, что лапласовское априорное распределение обычно приводит к более разреженным решающим правилам
- И регуляризованная линейная/логистическая регрессия, и метод релевантных векторов не инвариантны относительно линейных преобразований базисных функций

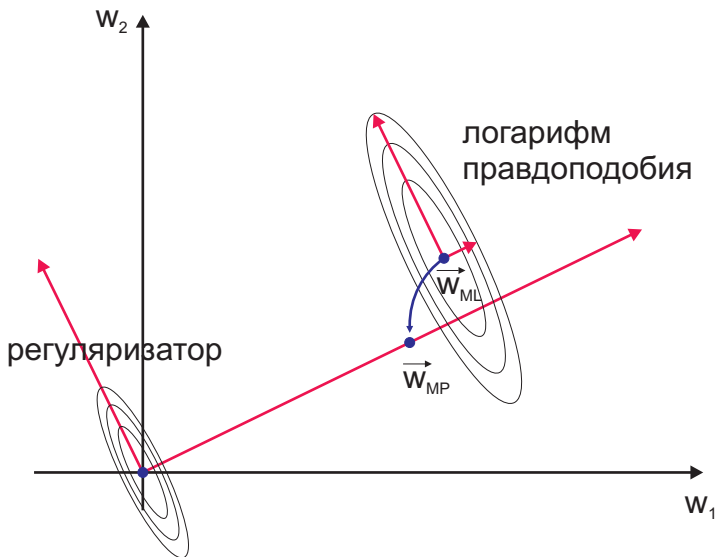
# Линейная неинвариантность

- Рассмотрим невырожденную матрицу  $L \in \mathbb{R}^{m \times m}$
- Пусть  $\psi(\mathbf{x}) = L\phi(\mathbf{x})$  — новое множество базисных функций
- Поскольку наш классификатор линеен по базисным функциям, вполне естественно требовать, чтобы классификатор, обученный по базисным функциям  $\psi(\mathbf{x})$ , был эквивалентен классификатору, полученному при использовании базисных функций  $\phi(\mathbf{x})$
- К сожалению, это не так в случае RVM и регуляризованной логистической регрессии

# Регуляризация в методе релевантных векторов



# Недиагональная регуляризация





# Недиагональная регуляризация

$$p(\mathbf{t}|X, \mathbf{w}) \approx p(\mathbf{t}|X, \mathbf{w}_{ML}) \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{w}_{ML})^T H(\mathbf{w} - \mathbf{w}_{ML})\right)$$

$$\exists Q: Q^T = Q^{-1}, QHQ^T = \Lambda = \text{diag}\{h_1, \dots, h_K\}$$

Пусть  $\mathbf{u}$  — новые переменные, ассоциированные с собственными векторами гессиана логарифма правдоподобия:  $\mathbf{u} = Q\mathbf{w}$ ,  $\mathbf{w} = Q^T\mathbf{u}$ .

Правдоподобие:

$$p(\mathbf{t}|X, \mathbf{u}) = p(\mathbf{t}|X, \mathbf{u}_{ML}) \exp\left(-\frac{1}{2}(\mathbf{u} - \mathbf{u}_{ML})^T \Lambda(\mathbf{u} - \mathbf{u}_{ML})\right)$$

Регуляризатор:

$$p(\mathbf{u}|\boldsymbol{\alpha}) = \sqrt{\frac{\det(A)}{(2\pi)^K}} \exp\left(-\frac{1}{2}\mathbf{u}^T A\mathbf{u}\right), A = \text{diag}(\alpha_1, \dots, \alpha_K), \alpha_j \geq 0.$$

# Сепарабельные функции в обоснованности

Обоснованность:

$$\int p(\mathbf{t}|X, \mathbf{u})p(\mathbf{u}|\boldsymbol{\alpha})d\mathbf{u} \rightarrow \max_{\boldsymbol{\alpha}}$$

Сепарабельность правдоподобия:

$$\begin{aligned} p(\mathbf{t}|X, \mathbf{u}) &= p(\mathbf{t}|X, \mathbf{u}_{ML}) \exp\left(-\frac{1}{2}(\mathbf{u} - \mathbf{u}_{ML})^T \Lambda (\mathbf{u} - \mathbf{u}_{ML})\right) = \\ p(\mathbf{t}|X, \mathbf{u}_{ML}) \exp\left(-\frac{1}{2} \sum_{j=1}^K h_j (u_j - u_{ML,j})^2\right) &= p(\mathbf{t}|X, \mathbf{u}_{ML}) \prod_{j=1}^K \underbrace{\exp\left(-\frac{1}{2} h_j (u_j - u_{ML,j})^2\right)}_{f(u_j)} \end{aligned}$$

Сепарабельность априорного распределения:  $p(\mathbf{u}|\boldsymbol{\alpha}) = \prod_{j=1}^K p(u_j|\alpha_j)$

$$\text{Обоснованность: } \int p(\mathbf{t}|X, \mathbf{u})p(\mathbf{u}|\boldsymbol{\alpha})d\mathbf{u} = \prod_{j=1}^K \underbrace{\int f(u_j)p(u_j|\alpha_j)du_j}_{g(\alpha_j)}$$

# Гауссовское априорное распределение

$$p(u_j|\alpha_j) = \mathcal{N}(u_j|0, \alpha_j^{-1})$$

$$g^G(\alpha_j) = \sqrt{\frac{\alpha_j}{2\pi}} \int \exp\left(-\frac{h_j}{2}(u_j - u_{ML,j})^2 - \frac{\alpha_j}{2}u_j^2\right) du_j = \\ \sqrt{\frac{\alpha_j}{h_j + \alpha_j}} \exp\left(-\frac{h_j\alpha_j u_{ML,j}^2}{2(h_j + \alpha_j)}\right),$$

$$\alpha_j^* = \begin{cases} \frac{h_j}{h_j u_{ML,j}^2 - 1}, & \text{если } h_j u_{ML,j}^2 > 1 \\ +\infty, & \text{иначе} \end{cases}$$

В данном случае условие на релевантность степени свободы можно получить в явном виде.

# Метод релевантных собственных векторов

**Вход:** Обучающая выборка  $\{\mathbf{x}_i, t_i\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $t_i \in \{+1, -1\}$ ;

Матрица обобщенных признаков  $\Phi = \{\phi_j(\mathbf{x}_i)\}_{i,j=1}^{n,m}$ ;

**Выход:** Набор весов  $\mathbf{w}_{MP}$  для решающего правила

$$t_*(\mathbf{x}) = \text{sign} \left( \sum_{j=1}^m w_{MP,j} \phi_j(\mathbf{x}) \right);$$

- 1: Найти  $\mathbf{w}_{ML} = \arg \max p(t|X, \mathbf{w})$ ;
- 2: Вычислить гессиан  $H = \nabla \nabla \log p(t|X, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_{ML}}$ ;
- 3: Вычислить собственные вектора и собственные значения гессиана  $-H = Q^T \Lambda Q$ ,  $\Lambda = \text{diag}(h_1, \dots, h_m)$ ;
- 4: Вычислить  $\mathbf{u}_{ML} = Q\mathbf{w}_{ML}$ ;
- 5: **для**  $j = 1, \dots, m$
- 6:     **если**  $h_j u_{ML,j}^2 > 1$  **то**
- 7:          $\alpha_j^* := \frac{h_j}{h_j u_{ML,j}^2 - 1}$ ;
- 8:     **иначе**
- 9:          $\alpha_j^* := +\infty$
- 10: Найти  $\mathbf{w}_{MP} = \arg \max p(t|X, \mathbf{w})p(Q\mathbf{w}|\boldsymbol{\alpha}^*)$

# Лапласовское априорное распределение

$$p(u_j|\alpha_j) = \mathcal{L}(u_j|\alpha_j^{-1})$$

Интеграл обоснованности может быть вычислен аналитически, разобьем интеграл на два

$$\begin{aligned} g^L(\alpha_j) &= \frac{\alpha_j}{4} \int_{-\infty}^0 \exp\left(-\frac{h_j(u_j - u_{ML,j})^2}{2} + \frac{\alpha_j}{2}u_j\right) du_j + \\ &\quad \frac{\alpha_j}{4} \int_0^{+\infty} \exp\left(-\frac{h_j(u_j - u_{ML,j})^2}{2} - \frac{\alpha_j}{2}u_j\right) du_j = \\ &\quad \frac{\alpha_j}{4} \sqrt{\frac{\pi}{2h_j}} \exp\left(-\frac{h_i u_{ML,i}^2}{2}\right) \times [\exp(x_1^2) \operatorname{erfc}(x_1) + \exp(x_2^2) \operatorname{erfc}(x_2)] \end{aligned}$$

Здесь  $x_1 = \sqrt{\frac{h_i}{2}} \left(\frac{\alpha_i}{2h_i} - u_{ML,i}\right)$ ,  $x_2 = \sqrt{\frac{h_i}{2}} \left(\frac{\alpha_i}{2h_i} + u_{ML,i}\right)$ ,

$$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^{+\infty} \exp(-\xi^2) d\xi$$

## Особенности численного вычисления

$$g^L(\alpha_j) = \frac{\alpha_j}{4} \sqrt{\frac{\pi}{2h_j}} \exp\left(-\frac{h_i u_{ML,i}^2}{2}\right) \times [\exp(x_1^2) \operatorname{erfc}(x_1) + \exp(x_2^2) \operatorname{erfc}(x_2)]$$

При  $x_{1,2}$  существенно отличных от нуля, возникают численные трудности с вычислением этого выражения. При  $x > 27$  выражение  $\exp(x^2) > 10^{300}$ , при  $x > 26$  выражение  $\operatorname{erfc}(x) < 10^{-300}$ .

Пусть  $x_j \gg 0$ ,  $j \in \{1, 2\}$ , тогда можно показать, что

$$\exp(x_j^2) \operatorname{erfc}(x_j) \approx 1/(\sqrt{\pi}x_j)$$

При  $x_j \ll 0$  объединяем  $\exp(-h_i u_{ML,i}^2/2)$  и  $\exp(x_j^2)$

$$\exp(-h_i u_{ML,i}^2/2) \exp(x_j^2) = \exp(y_j),$$

где

$$y_{1,2} = \frac{\alpha_i^2}{8h_i} \mp \frac{\alpha_i u_{ML,i}}{2}.$$

# Метод релевантных собственных векторов с лапласовским регуляризатором

**Вход:** Обучающая выборка  $\{\mathbf{x}_i, t_i\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $t_i \in \{+1, -1\}$ ;

Матрица обобщенных признаков  $\Phi = \{\phi_j(\mathbf{x}_i)\}_{i,j=1}^{n,m}$ ;

**Выход:** Набор весов  $\mathbf{w}_{MP}$  для решающего правила

$$t_*(\mathbf{x}) = \text{sign} \left( \sum_{j=1}^m w_{MP,j} \phi_j(\mathbf{x}) \right);$$

- 1: Найти  $\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(\mathbf{t}|X, \mathbf{w})$ ;
- 2: Вычислить гессиан  $H = \nabla \nabla \log p(\mathbf{t}|X, \mathbf{w})|_{\mathbf{w}=\mathbf{w}_{ML}}$ ;
- 3: Вычислить собственные вектора и собственные значения гессиана  $-H = Q^T \Lambda Q$ ,  $\Lambda = \text{diag}(h_1, \dots, h_m)$ ;
- 4: Вычислить  $\mathbf{u}_{ML} = Q \mathbf{w}_{ML}$ ;
- 5: для  $j = 1, \dots, m$
- 6:  $\alpha_j^* := \arg \max_{\alpha} f_j^L(h_j, u_{ML,j}, \alpha)$ ;
- 7: Найти  $\mathbf{u}_{MP} = \arg \max_{\mathbf{u}} p(\mathbf{t}|X, \mathbf{w}) p(\mathbf{u}|\boldsymbol{\alpha}^*)$ ;
- 8: Найти  $\mathbf{w}_{MP} = Q^T \mathbf{u}_{MP}$

# Оптимизация функции $g^L(\alpha_j)$

Точка максимума функции  $g^L(\alpha_j)$  не может быть выписана в явном виде, поэтому необходима численная оптимизация (впрочем, не слишком обременительная)

