

Разбор летучки

Лекция 10

Линейная регрессия

Екатерина Тузова

X — объекты в \mathbb{R}^n ; Y — ответы в \mathbb{R}

$X^l = (x_i, y_i)_{i=1}^l$ — обучающая выборка

$y_i = y(x_i)$, $y : X \rightarrow Y$ — неизвестная зависимость

X – объекты в \mathbb{R}^n ; Y – ответы в \mathbb{R}

$X^l = (x_i, y_i)_{i=1}^l$ – обучающая выборка

$y_i = y(x_i)$, $y : X \rightarrow Y$ – неизвестная зависимость

$a(x) = f(x, w)$ – модель зависимости,

$w \in \mathbb{R}^p$ – вектор параметров модели.

Многомерная линейная регрессия

x^1, \dots, x^n – числовые признаки

Модель многомерной линейной регрессии:

$$f(x, w) = \sum_{j=1}^n w_j x^j, \quad w \in \mathbb{R}$$

Многомерная линейная регрессия

x^1, \dots, x^n – числовые признаки

Модель многомерной линейной регрессии:

$$f(x, w) = \sum_{j=1}^n w_j x^j, \quad w \in \mathbb{R}$$

Функционал квадрата ошибки:

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 \rightarrow \min_w$$

Матричное представление

$$X_{l \times n} = \begin{pmatrix} x_1^1 & \dots & x_1^n \\ \dots & \dots & \dots \\ x_l^1 & \dots & x_l^n \end{pmatrix} \quad y_{l \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_l \end{pmatrix} \quad w_{n \times 1} = \begin{pmatrix} w_1 \\ \dots \\ w_n \end{pmatrix}$$

Матричное представление

$$X_{l \times n} = \begin{pmatrix} x_1^1 & \dots & x_1^n \\ \dots & \dots & \dots \\ x_l^1 & \dots & x_l^n \end{pmatrix} \quad y_{l \times 1} = \begin{pmatrix} y_1 \\ \dots \\ y_l \end{pmatrix} \quad w_{n \times 1} = \begin{pmatrix} w_1 \\ \dots \\ w_n \end{pmatrix}$$

Функционал квадрата ошибки:

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 = \|Xw - y\|^2 \rightarrow \min_w$$

Необходимое условие минимума:

$$\frac{\partial Q(w)}{\partial w} = 2X^T(Xw - y) = 0$$

Нормальная система уравнений

Необходимое условие минимума:

$$\frac{\partial Q(w)}{\partial w} = 2X^T(Xw - y) = 0$$

Откуда следует нормальная система задачи МНК:

$$X^T X w = X^T y$$

$X^T X$ – ковариационная матрица признаков x^1, \dots, x^n

Нормальная система уравнений

Нормальная система задачи МНК:

$$X^T X w = X^T y$$

Нормальная система уравнений

Нормальная система задачи МНК:

$$X^T X w = X^T y$$

Решение системы:

$$w^* = (X^T X)^{-1} X^T y = X^+ y$$

X^+ – псевдообратная матрица

Нормальная система уравнений

Нормальная система задачи МНК:

$$X^T X w = X^T y$$

Решение системы:

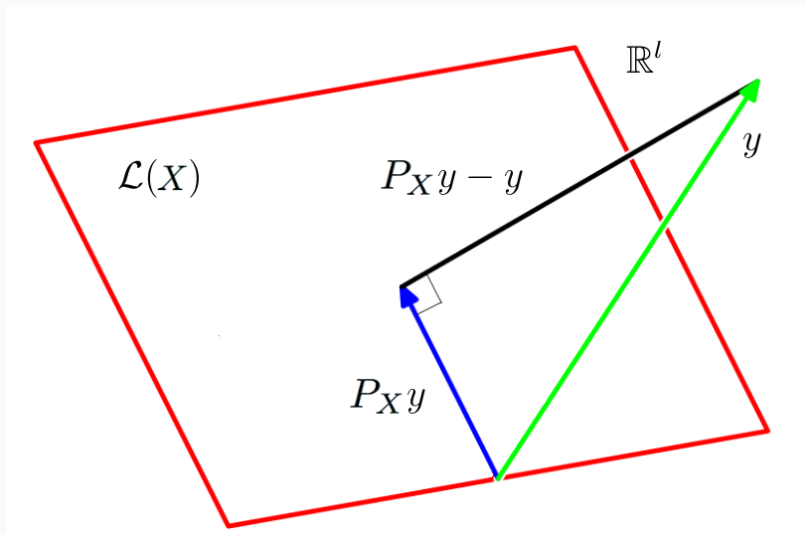
$$w^* = (X^T X)^{-1} X^T y = X^+ y$$

X^+ – псевдообратная матрица

Значение функционала: $Q(w^*) = \|P_X y - y\|^2$

где $P_X = X X^+ = X (X^T X)^{-1} X^T$ – проекционная матрица

Геометрический смысл



Сингулярное разложение (экономное)

Произвольная $l \times n$ -матрица представима в виде сингулярного разложения:

$$X = VDU^T$$

Сингулярное разложение (экономное)

Произвольная $l \times n$ -матрица представима в виде сингулярного разложения:

$$X = VDU^T$$

Основные свойства сингулярного разложения:

- $V_{l \times n} = (v_1, \dots, v_n)$ ортогональна, $V^T V = I_n$,
столбцы v_j - собственные векторы матрицы XX^T
- $U_{n \times n} = (u_1, \dots, u_n)$ ортогональна, $U^T U = I_n$,
столбцы u_j - собственные векторы матрицы $X^T X$
- D диагональна, $D_{n \times n} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$,
 $\lambda_j > 0$ - собственные значения матриц $X^T X$ и XX^T

Решение МНК через сингулярное разложение

Псевдообратная X^+ , вектор МНК-решения w^* , МНК-аппроксимация целевого вектора Xw^*

Псевдообратная X^+ , вектор МНК-решения w^* , МНК-аппроксимация целевого вектора Xw^*

$$X^+ = (UDV^TVDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T$$

Псевдообратная X^+ , вектор МНК-решения w^* , МНК-аппроксимация целевого вектора Xw^*

$$X^+ = (UDV^TVDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T$$
$$w^* = X^+y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y)$$

Решение МНК через сингулярное разложение

Псевдообратная X^+ , вектор МНК-решения w^* , МНК-аппроксимация целевого вектора Xw^*

$$X^+ = (UDV^TVDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T$$

$$w^* = X^+y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y)$$

$$Xw^* = P_X y = (VDU^T)UD^{-1}V^T y = VV^T y = \sum_{j=1}^n v_j (v_j^T y)$$

Решение МНК через сингулярное разложение

Псевдообратная X^+ , вектор МНК-решения w^* , МНК-аппроксимация целевого вектора Xw^*

$$X^+ = (UDV^TVDU^T)^{-1}UDV^T = UD^{-1}V^T = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j v_j^T$$

$$w^* = X^+y = UD^{-1}V^T y = \sum_{j=1}^n \frac{1}{\sqrt{\lambda_j}} u_j (v_j^T y)$$

$$Xw^* = P_X y = (VDU^T)UD^{-1}V^T y = VV^T y = \sum_{j=1}^n v_j (v_j^T y)$$

$$\|w^*\|^2 = \|UD^{-1}V^T y\|^2 = \sum_{j=1}^n \frac{1}{\lambda_j} (v_j^T y)^2$$

Если матрица $X^T X$ плохо обусловлена, то:

- решение становится неустойчивым и неинтерпретируемым, $\|w^*\|$ велико
- $E_{in} = \|Xw^* - y\|$ - мало
- $E_{out} = \|X'w^* - y'\|$ - велико

Если матрица $X^T X$ плохо обусловлена, то:

- решение становится неустойчивым и неинтерпретируемым, $\|w^*\|$ велико
- $E_{in} = \|Xw^* - y\|$ – мало
- $E_{out} = \|X'w^* - y'\|$ – велико

Вопрос: Как бороться с этой проблемой?

Стратегии устранения мультиколлинеарности и переобучения:

- Регуляризация: $\|w\| \rightarrow \min$
- Отбор признаков: $x^1, \dots, x^n \rightarrow x^{j_1}, \dots, x^{j_m}, \quad m \ll n$
- Преобразование признаков: $x^1, \dots, x^n \rightarrow g^1, \dots, g^m, \quad m \ll n$

Штраф за увеличение нормы вектора весов $\|w\|$:

$$Q_\tau(w) = \|Xw - y\|^2 + \tau\|w\|^2$$

где τ – неотрицательный параметр регуляризации.

Штраф за увеличение нормы вектора весов $\|w\|$:

$$Q_\tau(w) = \|Xw - y\|^2 + \tau\|w\|^2$$

где τ – неотрицательный параметр регуляризации.

Модифицированное МНК-решение (τI_n – «гребень»)

$$w_\tau^* = (X^T X + \tau I_n)^{-1} X^T y$$

Штраф за увеличение нормы вектора весов $\|w\|$:

$$Q_\tau(w) = \|Xw - y\|^2 + \tau\|w\|^2$$

где τ – неотрицательный параметр регуляризации.

Модифицированное МНК-решение (τI_n – «гребень»)

$$w_\tau^* = (X^T X + \tau I_n)^{-1} X^T y$$

Вопрос: Можно ли подбирать τ не вычисляя каждый раз обратную матрицу?

Модифицированное МНК-решение (τI_n — «гребень»)

$$w_{\tau}^* = (X^T X + \tau I_n)^{-1} X^T y$$

Модифицированное МНК-решение (τI_n — «гребень»)

$$w_{\tau}^* = (X^T X + \tau I_n)^{-1} X^T y$$

Преимущество сингулярного разложения:

Можно подбирать параметр τ , вычислив сингулярное разложение только один раз.

$$w_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y)$$

$$w_{\tau}^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y)$$

$$X w_{\tau}^* = V D U^T w_{\tau}^* = V \operatorname{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y)$$

$$w_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y)$$

$$X w_\tau^* = V D U^T w_\tau^* = V \operatorname{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y)$$

$$\|w_\tau^*\|^2 = \|U(D^2 + \tau I_n)^{-1} D V^T y\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2$$

$$w_\tau^* = U(D^2 + \tau I_n)^{-1} D V^T y = \sum_{j=1}^n \frac{\sqrt{\lambda_j}}{\lambda_j + \tau} u_j (v_j^T y)$$

$$X w_\tau^* = V D U^T w_\tau^* = V \operatorname{diag} \left(\frac{\lambda_j}{\lambda_j + \tau} \right) V^T y = \sum_{j=1}^n \frac{\lambda_j}{\lambda_j + \tau} v_j (v_j^T y)$$

$$\|w_\tau^*\|^2 = \|U(D^2 + \tau I_n)^{-1} D V^T y\|^2 = \sum_{j=1}^n \frac{\lambda_j}{(\lambda_j + \tau)^2} (v_j^T y)^2$$

$X w_\tau^* \neq X w^*$ – зато решение становится более устойчивым

Контрольная выборка: $X^k = (x'_i, y'_i)_{i=1}^k$

$$Q(w_\tau^*, X^k) = \|X'w_\tau^* - y'\|^2 = \|X'U \operatorname{diag}\left(\frac{\lambda_j}{\lambda_j + \tau}\right) V^T y - y'\|^2$$

Выбор параметра регуляризации τ

Контрольная выборка: $X^k = (x'_i, y'_i)_{i=1}^k$

$$Q(w_\tau^*, X^k) = \|X'w_\tau^* - y'\|^2 = \|X'U \operatorname{diag}\left(\frac{\lambda_j}{\lambda_j + \tau}\right) V^T y - y'\|^2$$

Зависимость $Q(\tau)$ обычно имеет характерный минимум.

$$\begin{cases} Q(w, X^l) = \|Xw - y\|^2 \rightarrow \min_w \\ \sum_{j=1}^n |w_j| \leq \kappa \end{cases}$$

$$\begin{cases} Q(w, X^l) = \|Xw - y\|^2 \rightarrow \min_w \\ \sum_{j=1}^n |w_j| \leq \kappa \end{cases}$$

После замены переменных:

$$\begin{cases} w_j = w_j^+ - w_j^- \\ |w_j| = w_j^+ + w_j^-, \quad w_j^+, w_j^- \geq 0 \end{cases}$$

ограничения принимают канонический вид:

$$\sum_{j=1}^n w_j^+ + w_j^- \leq \kappa$$

$$\begin{cases} Q(w, X^l) = \|Xw - y\|^2 \rightarrow \min_w \\ \sum_{j=1}^n |w_j| \leq \kappa \end{cases}$$

После замены переменных:

$$\begin{cases} w_j = w_j^+ - w_j^- \\ |w_j| = w_j^+ + w_j^-, \quad w_j^+, w_j^- \geq 0 \end{cases}$$

ограничения принимают канонический вид:

$$\sum_{j=1}^n w_j^+ + w_j^- \leq \kappa$$

Чем меньше κ , тем больше j таких, что $w_j^+ = w_j^- = 0$

Нелинейная регрессия

Вопрос: Что изменится, если модель регрессии не линейна?

$$f(x, w), \quad w \in \mathbb{R}^p$$

Начальное приближение $w^{(0)} = (w_1^{(0)}, \dots, w_p^{(0)})$

Итерационный процесс: $w^{(t+1)} = w^{(t)} - h_t(Q''(w^{(t)}))^{-1}Q'(w^{(t)})$

$Q'(w^{(t)})$ – градиент функционала Q в точке $w^{(t)}$

$Q''(w^{(t)})$ – гессиан функционала Q в точке $w^{(t)}$

h_t – величина шага

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 \rightarrow \min_w$$

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 \rightarrow \min_w$$

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial (f(x_i, w))}{\partial w_j}$$

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 \rightarrow \min_w$$

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial (f(x_i, w))}{\partial w_j}$$

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = 2 \sum_{i=1}^l \frac{\partial (f(x_i, w))}{\partial w_j} \frac{\partial (f(x_i, w))}{\partial w_k} - 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial^2 (f(x_i, w))}{\partial w_j \partial w_k}$$

$$Q(w, X^l) = \sum_{i=1}^l (f(x_i, w) - y_i)^2 \rightarrow \min_w$$

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial(f(x_i, w))}{\partial w_j}$$

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = 2 \sum_{i=1}^l \frac{\partial(f(x_i, w))}{\partial w_j} \frac{\partial(f(x_i, w))}{\partial w_k} - 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial^2(f(x_i, w))}{\partial w_j \partial w_k}$$

Вопрос: Какая часть самая тяжелая?

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial (f(x_i, w))}{\partial w_j}$$

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial(f(x_i, w))}{\partial w_j}$$

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = 2 \sum_{i=1}^l \frac{\partial(f(x_i, w))}{\partial w_j} \frac{\partial(f(x_i, w))}{\partial w_k} - 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial^2(f(x_i, w))}{\partial w_j \partial w_k}$$

$$\frac{\partial Q(w)}{\partial w_j} = 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial(f(x_i, w))}{\partial w_j}$$

$$\frac{\partial^2 Q(w)}{\partial w_j \partial w_k} = 2 \sum_{i=1}^l \frac{\partial(f(x_i, w))}{\partial w_j} \frac{\partial(f(x_i, w))}{\partial w_k} - 2 \sum_{i=1}^l (f(x_i, w) - y_i) \frac{\partial^2(f(x_i, w))}{\partial w_j \partial w_k}$$

Линеаризация $f(x_i, w)$ в окрестности текущего $w^{(t)}$:

$$f(x_i, w) = f(x_i, w^{(t)}) + \sum_{j=1}^p \frac{\partial(f(x_i, w))}{\partial w_j} (w_j - w_j^{(t)}) + o(w_j - w_j^{(t)})$$

⇒ второе слагаемое в гессиане обнулилось

Матричные обозначения

$$X_t = \left(\frac{\partial(f(x_i, w^{(t)}))}{\partial w_j^{(t)}} \right)_{l \times p} \quad \text{– матрица первых производных}$$
$$f_t = (f(x_i, w^{(t)}))_{l \times 1} \quad \text{– вектор значений } f$$

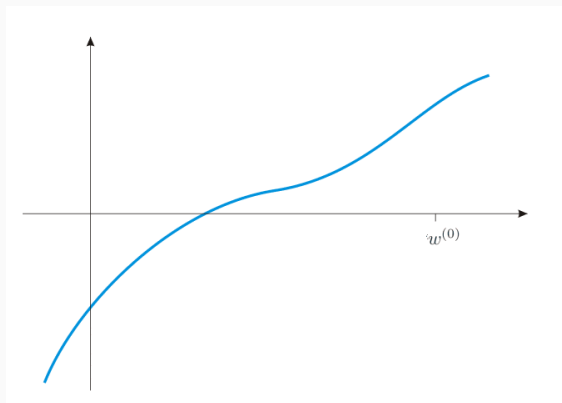
Формула t -й итерации метода Ньютона–Гаусса:

$$w^{(t+1)} = w^{(t)} - h_t \underbrace{(X_t^T X_t)^{-1} X_t^T (f_t - y)}_{\beta}$$

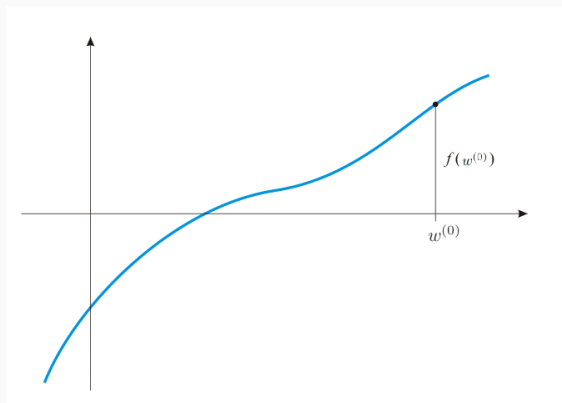
где β – решение многомерной линейной регрессии

$$\|X_t \beta - (f_t - y)\|^2 \rightarrow \min_{\beta}$$

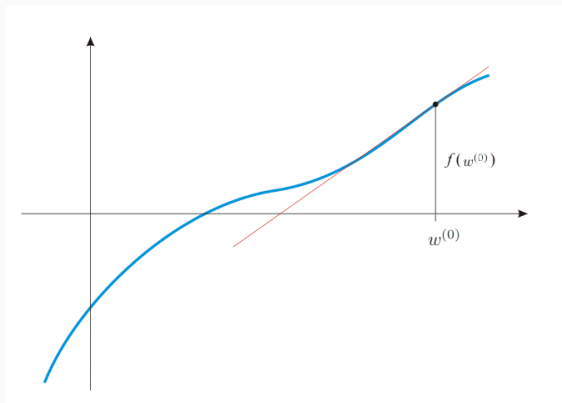
Метод Ньютона-Рафсена



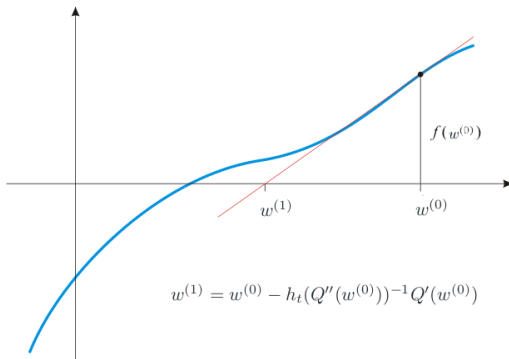
Метод Ньютона-Рафсена



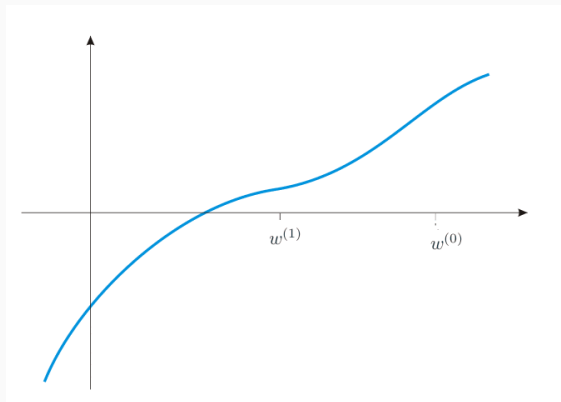
Метод Ньютона-Рафсена



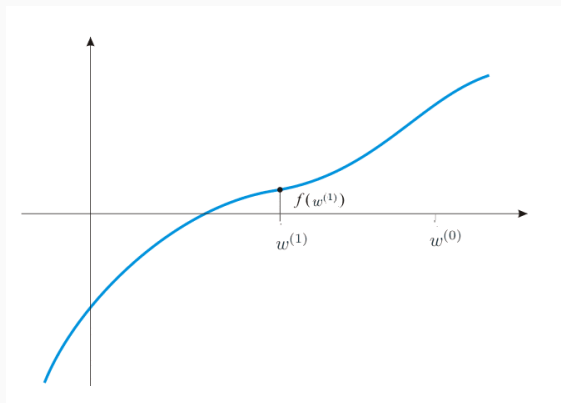
Метод Ньютона-Рафсена



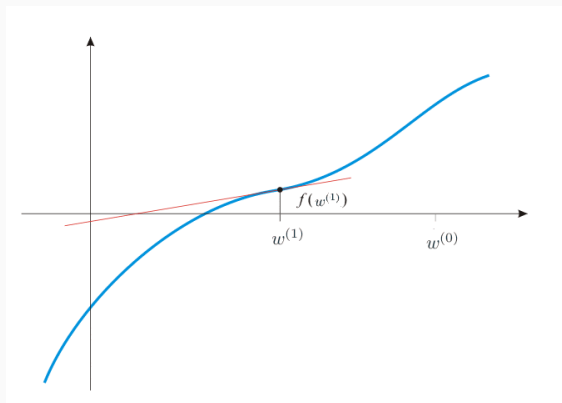
Метод Ньютона-Рафсена



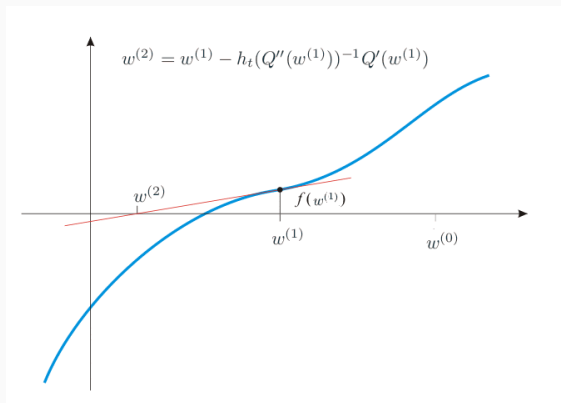
Метод Ньютона-Рафсена



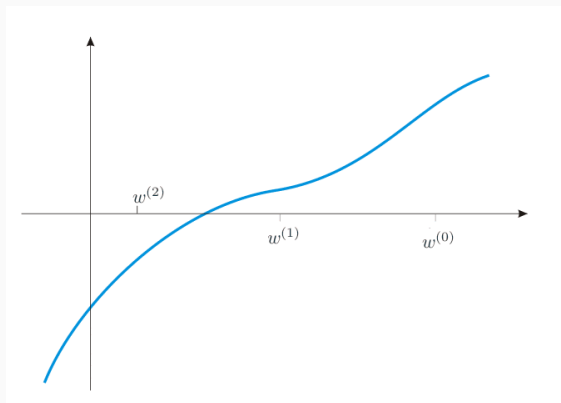
Метод Ньютона-Рафсена



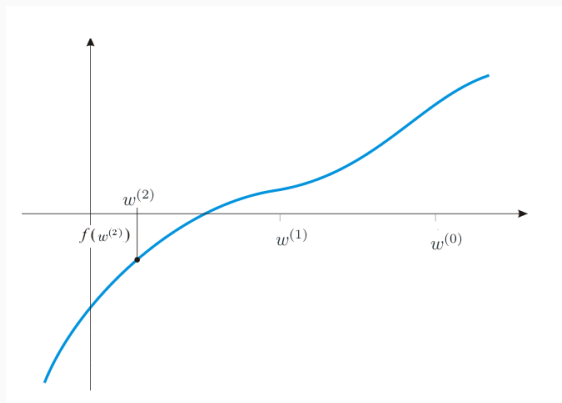
Метод Ньютона-Рафсена



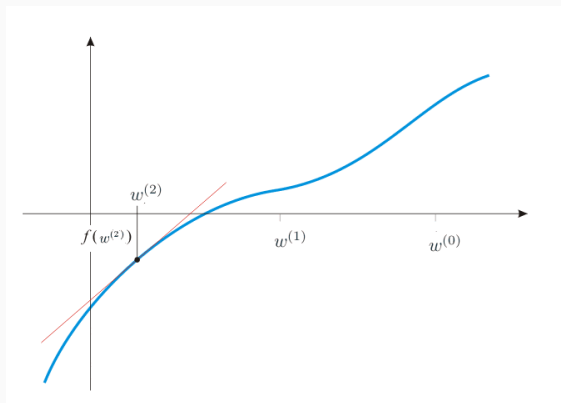
Метод Ньютона-Рафсена



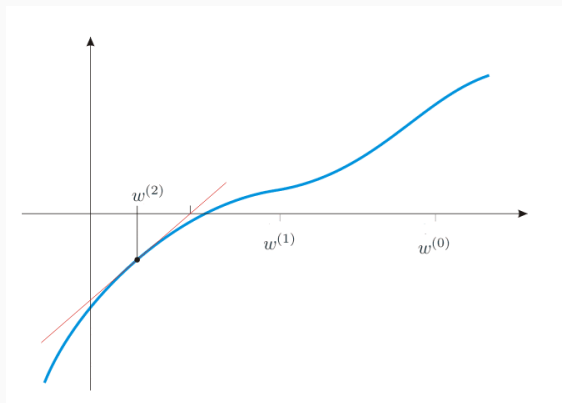
Метод Ньютона-Рафсена



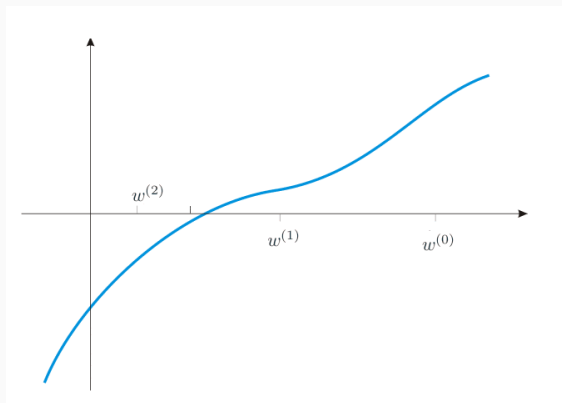
Метод Ньютона-Рафсена



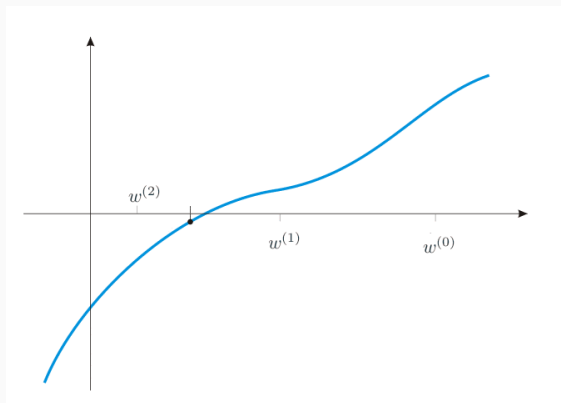
Метод Ньютона-Рафсена



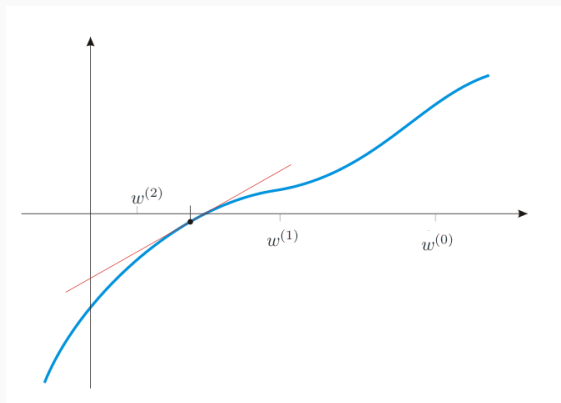
Метод Ньютона-Рафсена



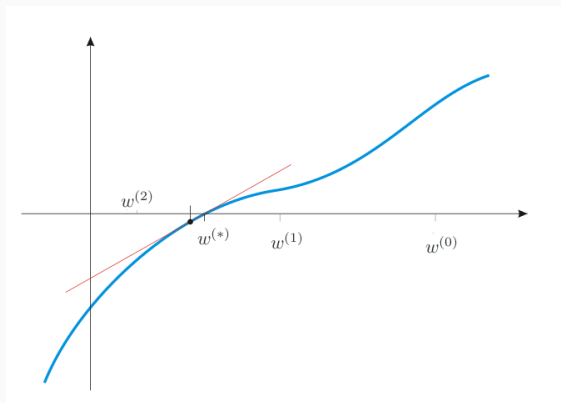
Метод Ньютона-Рафсена



Метод Ньютона-Рафсена



Метод Ньютона-Рафсена



Вопросы?

Что почитать по этой лекции

- T. Hastie, R. Tibshirani "The Elements of Statistical Learning"
Chapter 3

На следующей лекции

- Bias-variance tradeoff
- Кривые обучения