

Optymalizacja ciągła

5. Metoda stochastycznego spadku wzdłuż gradientu

Wojciech Kotłowski

Instytut Informatyki PP

<http://www.cs.put.poznan.pl/wkotlowski/>

04.04.2019

Wprowadzenie

Minimalizacja **różniczkowalnej** funkcji $f(\mathbf{w})$ **bez ograniczeń**:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w})$$

Uwaga: na potrzeby tego wykładu zmieniamy oznaczenie wektora zmiennych z \mathbf{x} na \mathbf{w}

W wielu zastosowaniach optymalizacji funkcja celu ma postać sumy po N elementach:

$$f(\mathbf{w}) = \sum_{i=1}^N f_i(\mathbf{w}),$$

To założenie będzie nam towarzyszyło przez cały wykład.

Przykład: regresja liniowa

Przewidywania/wyjaśnienie zmian ciągłej zmiennej wyjściowej (Y) pod wpływem zmian innych zmiennych (X_1, \dots, X_n).

Przykłady

- X – ceny akcji w ostatnim tygodniu, Y – cena akcji jutro.
- X – wyniki testów medycznych, Y – poziom zaawansowania choroby.
- X – wielkość programu, Y – czas pisania programu.
- X – warunki na drodze, czas, lokalizacja, Y – średnia prędkość samochodów.
- X – cechy domu Y – cena domu.

Przykład: regresja liniowa

Modelujemy zmienną Y jako funkcję liniową X_1, \dots, X_n :

$$Y = w_0 + w_1 X_1 + \dots + w_n X_n = \mathbf{X}^\top \mathbf{w}$$

gdzie $\mathbf{w} = (w_0, w_1, \dots, w_n)$ jest wektorem **współczynników** (zmiennych decyzyjnych), natomiast $\mathbf{X} = (1, X_1, \dots, X_n)$.

Przykład: regresja liniowa

Modelujemy zmienną Y jako funkcję liniową X_1, \dots, X_n :

$$Y = w_0 + w_1 X_1 + \dots + w_n X_n = \mathbf{X}^\top \mathbf{w}$$

gdzie $\mathbf{w} = (w_0, w_1, \dots, w_n)$ jest wektorem **współczynników** (zmiennych decyzyjnych), natomiast $\mathbf{X} = (1, X_1, \dots, X_n)$.

Wartości współczynników dopasujemy **na podstawie zbioru danych**

Przykład: regresja liniowa

Otrzymujemy zbiór danych historycznych, na którym znane są wartości y :

$$(x_{11}, x_{12}, \dots, x_{1n}, y_1)$$

$$(x_{21}, x_{22}, \dots, x_{2n}, y_2)$$

...

$$(x_{N1}, x_{N2}, \dots, x_{Nn}, y_N)$$

$$(\mathbf{x}_1, y_1)$$

$$(\mathbf{x}_2, y_2)$$

...

$$(\mathbf{x}_N, y_N)$$

lub w skrócie

Przykład: regresja liniowa

Otrzymujemy zbiór danych historycznych, na którym znane są wartości y :

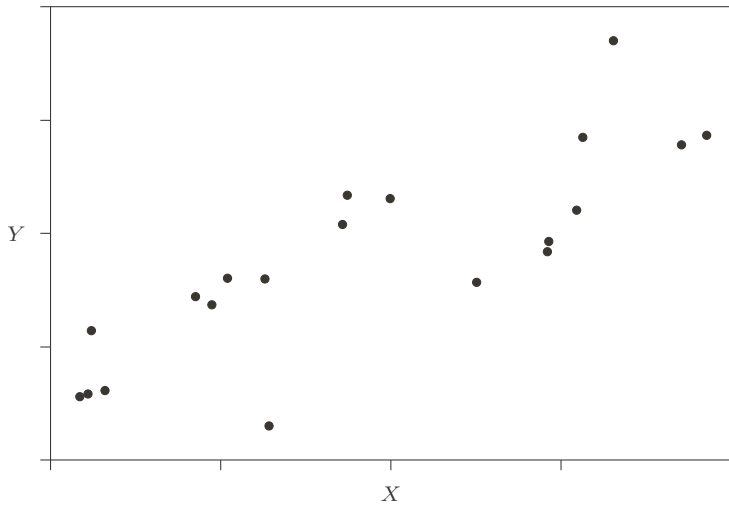
$$\begin{array}{ccc} (x_{11}, x_{12}, \dots, x_{1n}, y_1) & & (\mathbf{x}_1, y_1) \\ (x_{21}, x_{22}, \dots, x_{2n}, y_2) & & (\mathbf{x}_2, y_2) \\ \dots & \text{lub w skrócie} & \dots \\ (x_{N1}, x_{N2}, \dots, x_{Nn}, y_N) & & (\mathbf{x}_N, y_N) \end{array}$$

Wyznaczamy współczynniki \mathbf{w} tak, aby funkcja liniowa jak najlepiej przybliżała wartości zmiennej Y na danych, tzn. aby:

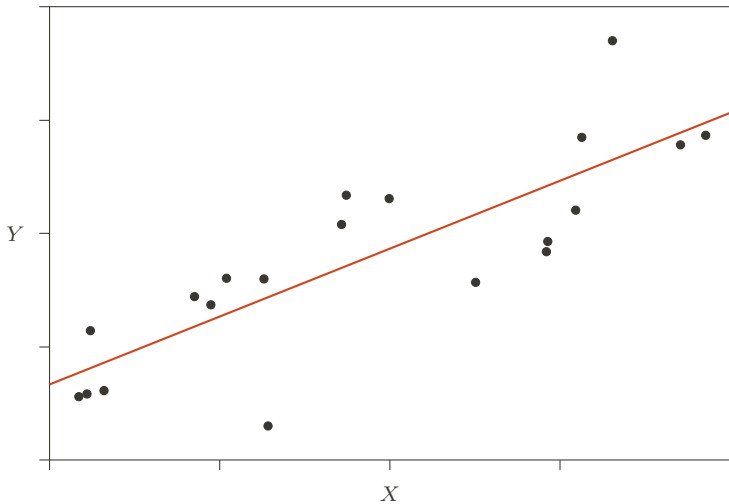
$$\hat{y}_i = \mathbf{x}_i^\top \mathbf{w}$$

było jak najbliższe y_i dla wszystkich $i = 1, \dots, N$

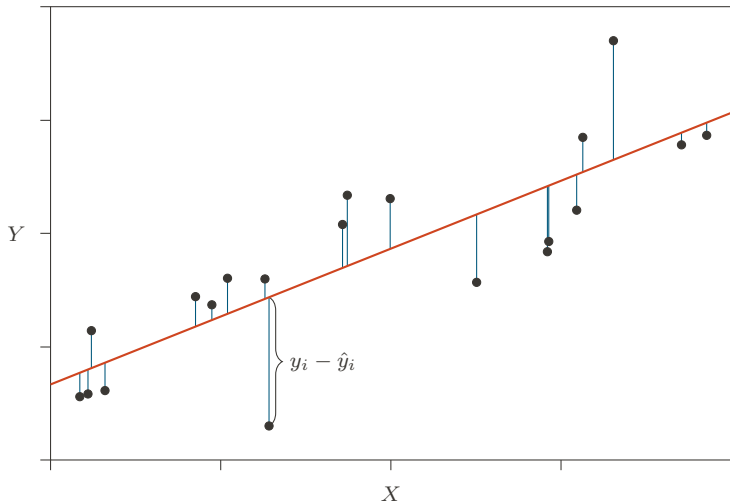
Przykład: regresja liniowa



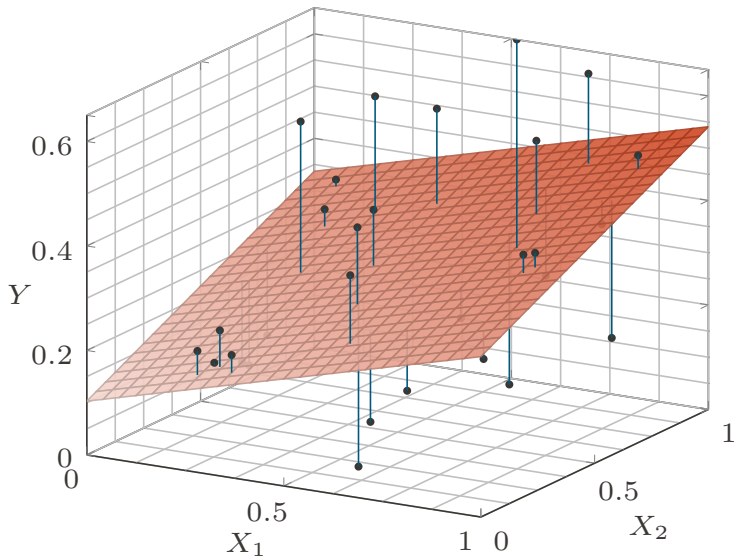
Przykład: regresja liniowa



Przykład: regresja liniowa



Przykład: regresja liniowa



Przykład: regresja liniowa

Metoda najmniejszych kwadratów:

Wyznacz współczynniki regresji w by minimalizować sumę kwadratów odchyłeń modelu od danych:

$$\min_{w \in \mathbb{R}^n} f(w) = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Przykład: regresja liniowa

Metoda najmniejszych kwadratów:

Wyznacz współczynniki regresji \mathbf{w} by minimalizować sumę kwadratów odchyłeń modelu od danych:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

Przykład: regresja liniowa

Metoda najmniejszych kwadratów:

Wyznacz współczynniki regresji \mathbf{w} by minimalizować sumę kwadratów odchyłeń modelu od danych:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Przykład: regresja liniowa

Metoda najmniejszych kwadratów:

Wyznacz współczynniki regresji \mathbf{w} by minimalizować sumę kwadratów odchyłeń modelu od danych:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Innymi słowy: wyznaczamy współczynniki modelu liniowego minimalizując **sumaryczny błąd** (wyrażony w postaci funkcji kwadratowej) **na całym zbiorze danych**

Przykłady: uczenie maszynowe

$$f(\mathbf{w}) = \sum_{i=1}^N f_i(\mathbf{w})$$

Zdecydowana większość metod **uczenia maszynowego** opiera się zasadniczo na tym samym schemacie:

Mając **zbiór danych**, wyznacz **parametry** modelu predykcyjnego **minimalizując sumaryczny błąd modelu na danych**

- \mathbf{w} – parametry modelu (współczynniki regresji, wagi w sieci neuronowej, itp.)
- $f_i(\mathbf{w})$ – błąd na pojedynczej obserwacji ze zbioru danych (np. błąd kwadratowy w regresji)
- $f(\mathbf{w})$ – sumaryczna wartość błędu na całym zbiorze danych

Podejście klasyczne

Znajdź rozwiązanie problemu:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w})$$

stosując jedną z klasycznych metod optymalizacji, np. metodę spadku wzdłuż gradientu lub Newtona-Raphsona

Podjęcie klasyczne

Znajdź rozwiązanie problemu:

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w})$$

stosując jedną z klasycznych metod optymalizacji, np. metodę spadku wzdłuż gradientu lub Newtona-Raphsona

Typowe problemy w praktyce:

- Policzenie gradientu/hesjanu wymaga sumowania N gradientów/hesjanów składowych funkcji celu:

$$\nabla f(\mathbf{w}) = \sum_{i=1}^N \nabla f_i(\mathbf{w})$$

- Wymiar problemu n może być bardzo duży, przez co metoda Newtona-Raphsona może być zbyt kosztowna do uruchomienia

Podójście klasyczne: spadek wzdłuż gradientu

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \underbrace{\sum_{i=1}^N \nabla f_i(\mathbf{w}_k)}_{\nabla f(\mathbf{w}_k)}$$

Wykonanie jednego kroku: przejście po wszystkich składowych funkcji celu.

Podójście klasyczne: spadek wzdłuż gradientu

Rozpoczynając od w_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

$$w_{k+1} = w_k - \alpha_k \underbrace{\sum_{i=1}^N \nabla f_i(w_k)}_{\nabla f(w_k)}$$

Wykonanie jednego kroku: przejście po wszystkich składowych funkcji celu.

Pomysł: policz gradient na **losowej pojedynczej składowej** f_i zamiast na całej funkcji f !

Metoda stochastycznego spadku wzdłuż gradientu (*stochastic gradient descent, SGD*)

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

Wylosuj jedną ze składowych $i_k \in \{1, \dots, N\}$

Wykonaj krok wzdłuż ujemnego gradientu funkcji $f_{i_k}(\mathbf{w}_k)$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k)$$

Metoda stochastycznego spadku wzdłuż gradientu (*stochastic gradient descent, SGD*)

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

Wylosuj jedną ze składowych $i_k \in \{1, \dots, N\}$

Wykonaj krok wzdłuż ujemnego gradientu funkcji $f_{i_k}(\mathbf{w}_k)$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k)$$

Dlaczego miałyby to w ogóle działać?

Metoda stochastycznego spadku wzdłuż gradientu (*stochastic gradient descent, SGD*)

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

Wylosuj jedną ze składowych $i_k \in \{1, \dots, N\}$

Wykonaj krok wzdłuż ujemnego gradientu funkcji $f_{i_k}(\mathbf{w}_k)$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k)$$

Dlaczego miałyby to w ogóle działać?

Możemy myśleć o gradiencie pojedynczej składowej $\nabla f_i(\mathbf{w}_k)$ jako o **oszacowaniu** całościowego gradientu $\nabla f(\mathbf{w}_k)$

Rzeczywiście, **średnio** gradient będzie wynosił:

$$\frac{1}{n} \nabla f_1(\mathbf{w}_k) + \frac{1}{n} \nabla f_2(\mathbf{w}_k) + \dots + \frac{1}{n} \nabla f_N(\mathbf{w}_k) = \frac{1}{n} \nabla f(\mathbf{w}_k)$$

Przykład: problem regresji liniowej

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Przykład: problem regresji liniowej

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Liczmy pochodne cząstkowe pojedynczej funkcji $f_i(\mathbf{w})$:

$$\frac{\partial f_i(\mathbf{w})}{\partial w_j} = -2(y_i - \mathbf{x}_i^\top \mathbf{w})x_{ij}$$

Przykład: problem regresji liniowej

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Liczmy pochodne cząstkowe pojedynczej funkcji $f_i(\mathbf{w})$:

$$\frac{\partial f_i(\mathbf{w})}{\partial w_j} = -2(y_i - \mathbf{x}_i^\top \mathbf{w})x_{ij}$$

Gradient wynosi więc: $\nabla f_i(\mathbf{w}) = -2(y_i - \mathbf{x}_i^\top \mathbf{w})\mathbf{x}_i$

Przykład: problem regresji liniowej

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Liczmy pochodne cząstkowe pojedynczej funkcji $f_i(\mathbf{w})$:

$$\frac{\partial f_i(\mathbf{w})}{\partial w_j} = -2(y_i - \mathbf{x}_i^\top \mathbf{w})x_{ij}$$

Gradient wynosi więc: $\nabla f_i(\mathbf{w}) = -2(y_i - \mathbf{x}_i^\top \mathbf{w})\mathbf{x}_i$

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

Wylosuj jedną ze obserwacji $i_k \in \{1, \dots, N\}$

Wykonaj krok wzdłuż ujemnego gradientu funkcji $f_{i_k}(\mathbf{w}_k)$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + 2\alpha_k(y_{i_k} - \mathbf{x}_{i_k}^\top \mathbf{w}_k)\mathbf{x}_{i_k}$$

Przykład: problem regresji liniowej

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \sum_{i=1}^N \underbrace{(y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{f_i(\mathbf{w})}$$

Liczmy pochodne cząstkowe pojedynczej funkcji $f_i(\mathbf{w})$:

$$\frac{\partial f_i(\mathbf{w})}{\partial w_j} = -2(y_i - \mathbf{x}_i^\top \mathbf{w})x_{ij}$$

Gradient wynosi więc: $\nabla f_i(\mathbf{w}) = -2(y_i - \mathbf{x}_i^\top \mathbf{w})\mathbf{x}_i$

Rozpoczynając od \mathbf{w}_1 , w kolejnych iteracjach $k = 1, 2, \dots$:

Wylosuj jedną ze obserwacji $i_k \in \{1, \dots, N\}$

Wykonaj krok wzdłuż ujemnego gradientu funkcji $f_{i_k}(\mathbf{w}_k)$:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + 2\alpha_k(y_{i_k} - \mathbf{x}_{i_k}^\top \mathbf{w}_k)\mathbf{x}_{i_k}$$

Krok w kierunku wektora obserwacji \mathbf{x}_{i_k} , z długością kroku proporcjonalną do „przeszacowania predykcji” $(y_{i_k} - \mathbf{x}_{i_k}^\top \mathbf{w}_k)$.

Zalety i wady algorytmu SGD

Zalety

- **Szybkość:** obliczenie gradientu wymaga wzięcia tylko jednej obserwacji.
- **Skalowalność:** cały zbiór danych nie musi nawet znajdować się w pamięci operacyjnej.
- **Prostota:** gradient funkcji f_i daje bardzo prosty wzór na modyfikację wag.

Wady

- **Wolna zbieżność:** algorytm w ogólności zbiega wolno i wymaga wielu iteracji po zbiorze uczącym.
- **Ustalenie długości kroku α_k :** wyznaczenie α_k przez jednowymiarową optymalizację nie przynosi dobrych rezultatów, ponieważ nie optymalizujemy oryginalnej funkcji f tylko jeden jej składnik f_i .

Zalety znacznie przewyższają wady: Algorytm SGD jest obecnie najczęściej używanym algorytmem w uczeniu maszynowym!

SGD w praktyce

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :
 - ▶ Ustalamy **stałą wartość** $\alpha_k = \alpha$

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :
 - ▶ Ustalamy **stałą wartość** $\alpha_k = \alpha$
⇒ Zwykle tak się robi w praktyce, działa dobrze ale wymaga ustalenia α metodą prób i błędów

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :
 - ▶ Ustalamy **stałą wartość** $\alpha_k = \alpha$
 \implies Zwykle tak się robi w praktyce, działa dobrze ale wymaga ustalenia α metodą prób i błędów
 - ▶ Bierzymy wartość kroku **malejącą jak** $\sim \frac{1}{\sqrt{k}}$: $\alpha_k = \alpha/\sqrt{k}$

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :
 - ▶ Ustalamy **stałą wartość** $\alpha_k = \alpha$
⇒ Zwykle tak się robi w praktyce, działa dobrze ale wymaga ustalenia α metodą prób i błędów
 - ▶ Bierzymy wartość kroku **malejącą jak** $\sim \frac{1}{\sqrt{k}}$: $\alpha_k = \alpha/\sqrt{k}$
⇒ Zapewniona zbieżność, ale czasem może zbiegać zbyt wolno.

SGD w praktyce

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności (działa bardzo podobnie, a daje znacznie prostszą implementację)
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku α_k :
 - ▶ Ustalamy **stałą wartość** $\alpha_k = \alpha$
⇒ Zwykle tak się robi w praktyce, działa dobrze ale wymaga ustalenia α metodą prób i błędów
 - ▶ Bierzymy wartość kroku **malejącą jak** $\sim \frac{1}{\sqrt{k}}$: $\alpha_k = \alpha/\sqrt{k}$
⇒ Zapewniona zbieżność, ale czasem może zbiegać zbyt wolno.
- W praktyce często liczy się gradient nie na pojedynczej obserwacji, ale na ich niewielkiej grupie (tzw. *mini-batch*)

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

wpływ „warunków początkowych”
(odległość od optimum na starcie)
(maleje z α)

wpływ wielkości gradientów
na trajektorii algorytmu
(rośnie z α)

algorytm stochastycznego spadku wzdłuż gradientu ze stałą długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$(\dots) \frac{1}{\alpha K} \simeq \frac{\alpha}{K} \sum_{k=1}^K (\dots)$$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$(\dots) \frac{1}{\alpha K} \simeq \alpha(\dots)$$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$(\dots) \frac{1}{K} \simeq \alpha^2$$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$\alpha \simeq \frac{(\dots)}{\sqrt{K}}$$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$\alpha \simeq \frac{(\dots)}{\sqrt{K}}$$

Stąd stosowane w praktyce malejące wartości α_k rzędu $\alpha_k = \frac{\text{const}}{\sqrt{k}}$

Metoda stochastycznego spadku wzdłuż gradientu

Twierdzenie: Niech $f_i(\mathbf{w})$ ($i = 1, \dots, N$) będą funkcjami wypukłymi i niech \mathbf{w}^* będzie minimum (globalnym) funkcji $f(\mathbf{w})$. Po K iteracjach algorytm stochastycznego spadku wzdłuż gradientu ze **stałą** długością kroku α gwarantuje:

$$\min_{k=1, \dots, K} f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{\|\mathbf{w}_1 - \mathbf{w}^*\|^2}{2\alpha K} + \frac{\alpha}{2K} \sum_{k=1}^K \|\nabla f_{i_k}(\mathbf{w}_k)\|^2$$

Wniosek: sensowne jest wzięcie α tak, aby oba człony były podobnej wielkości:

$$\alpha \simeq \frac{(\dots)}{\sqrt{K}}$$

Stąd stosowane w praktyce malejące wartości α_k rzędu $\alpha_k = \frac{\text{const}}{\sqrt{k}}$

Powyższe wartości α_k dają gwarancje na suboptymalność rzędu $O\left(\frac{1}{\sqrt{K}}\right)$

Porównanie zbieżności

algorytm	zbieżność	błąd po k iteracjach
Metoda Cauchy'ego	liniowa	$\propto e^{-ck}$
Metoda Newtona-Raphsona	kwadratowa	$\propto e^{-c_1 e^{c_2 k}}$
SGD	subliniowa	$\propto \frac{c}{\sqrt{k}}$

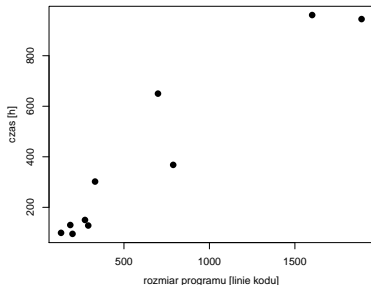
Algorytm SGD jest **bardzo wolny** w stosunku do uprzednio poznanych. Stosuje się go wtedy, gdy nie jest nam potrzebna duża dokładność punktu optimum, a zależy nam na czasie optymalizacji.

Przykład: szacowanie czasu pracy programistów

X	Y
Rozmiar programu	Oszacowany czas
186	130
699	650
132	99
272	150
291	128
331	302
199	95
1890	945
788	368
1601	961

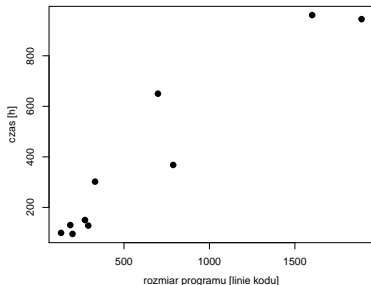
Przykład: szacowanie czasu pracy programistów

X	Y
Rozmiar programu	Oszacowany czas
186	130
699	650
132	99
272	150
291	128
331	302
199	95
1890	945
788	368
1601	961



Przykład: szacowanie czasu pracy programistów

X	Y
Rozmiar programu	Oszacowany czas
186	130
699	650
132	99
272	150
291	128
331	302
199	95
1890	945
788	368
1601	961



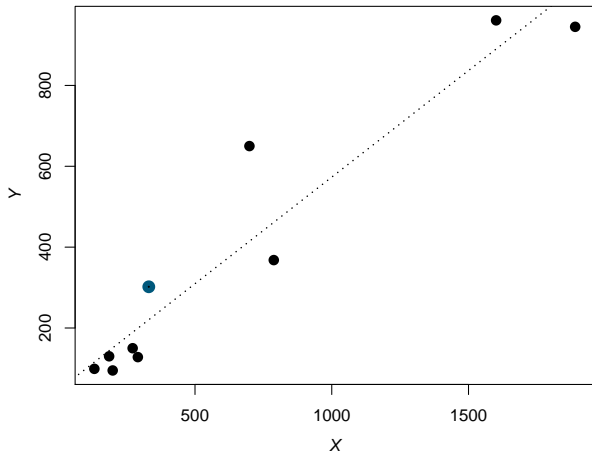
$$w_0 = 45.93, w_1 = 0.5273$$

Szacowanie czasu programistów

- Zaczynamy z rozwiązaniem $\mathbf{w} = (w_0, w_1) = \mathbf{0}$.
- Krok: $\alpha_k = 0.5/\sqrt{k}$.
- Jednostki zostały zamienione na 1000min i 1000 lini kodu, aby uniknąć dużych liczb.
(ujednolicenie skali jest w praktyce bardzo istotne dla zbieżności metody)

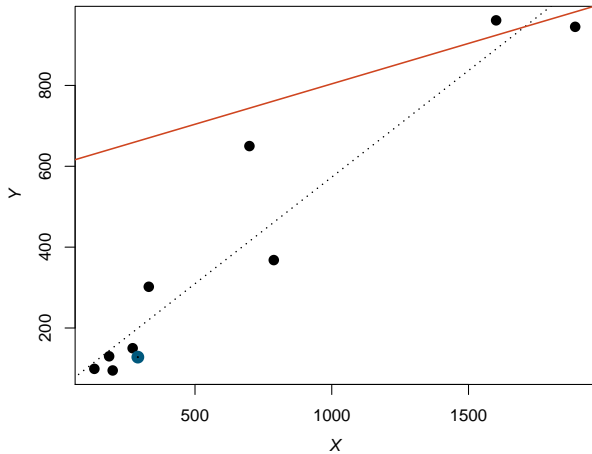
Szacowanie czasu programistów

iteracja 1



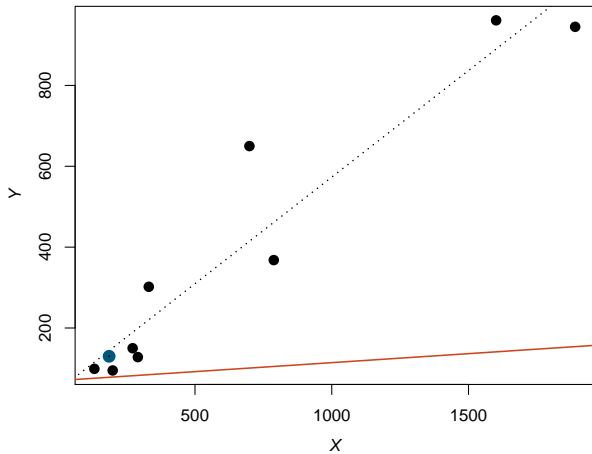
Szacowanie czasu programistów

iteracja 2



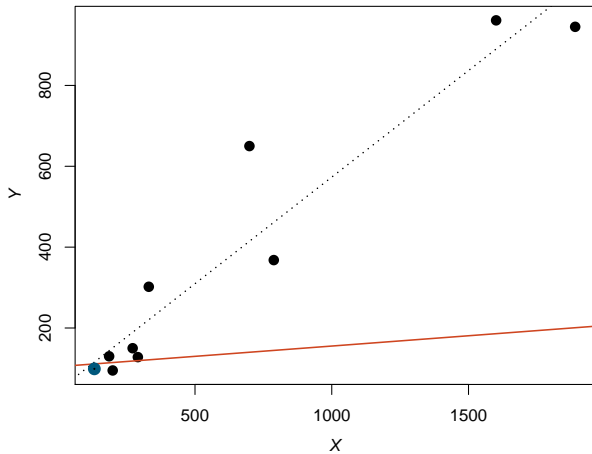
Szacowanie czasu programistów

iteracja 3



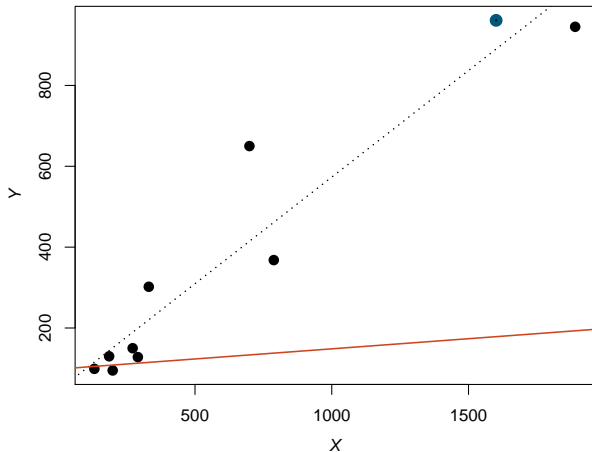
Szacowanie czasu programistów

iteracja 4



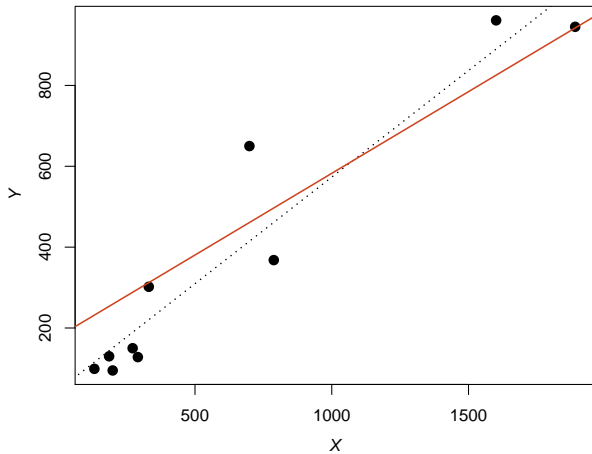
Szacowanie czasu programistów

iteracja 5



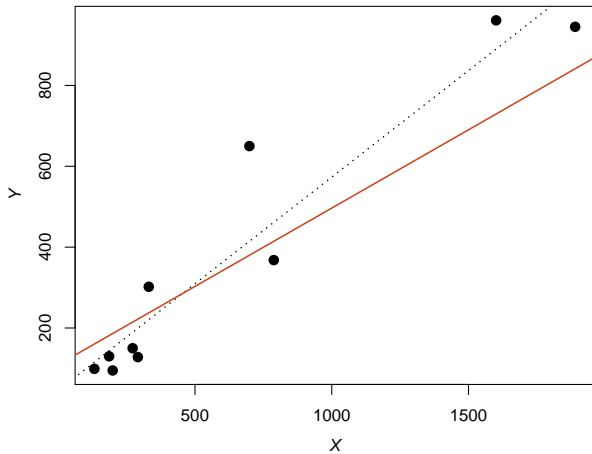
Szacowanie czasu programistów

iteracja 10



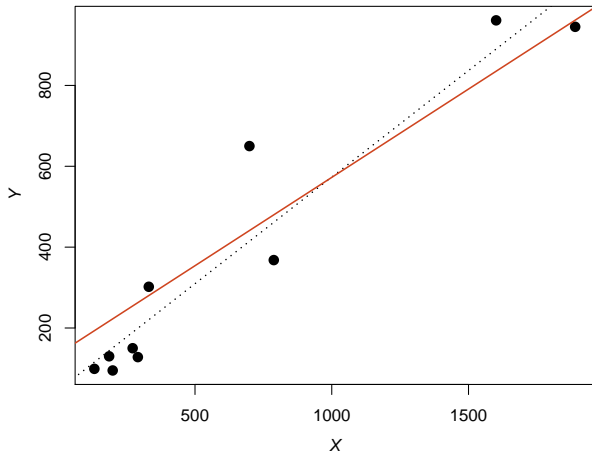
Szacowanie czasu programistów

iteracja 15



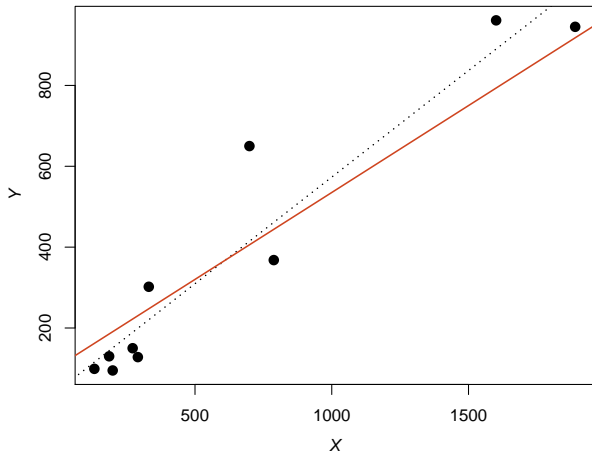
Szacowanie czasu programistów

iteracja 20



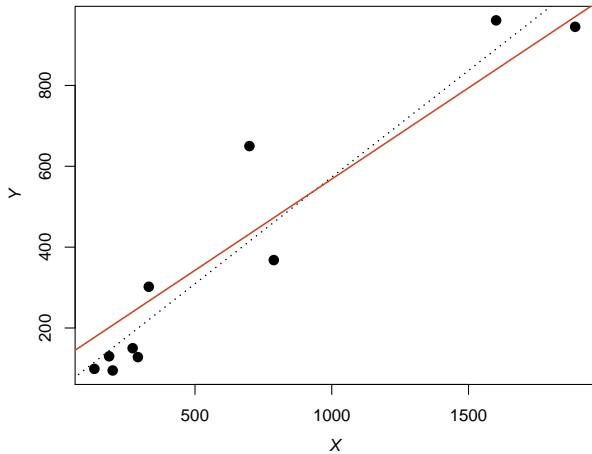
Szacowanie czasu programistów

iteracja 25



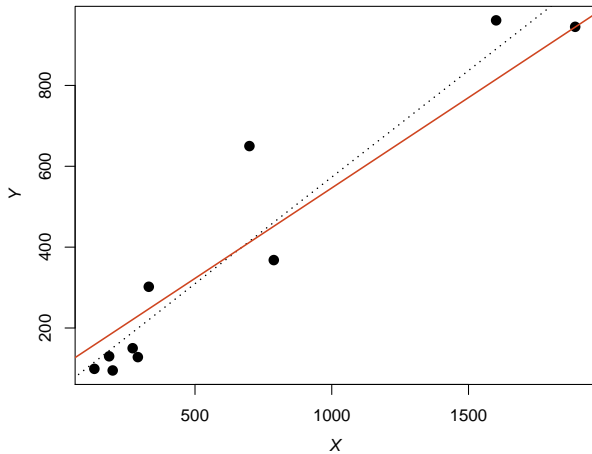
Szacowanie czasu programistów

iteracja 30



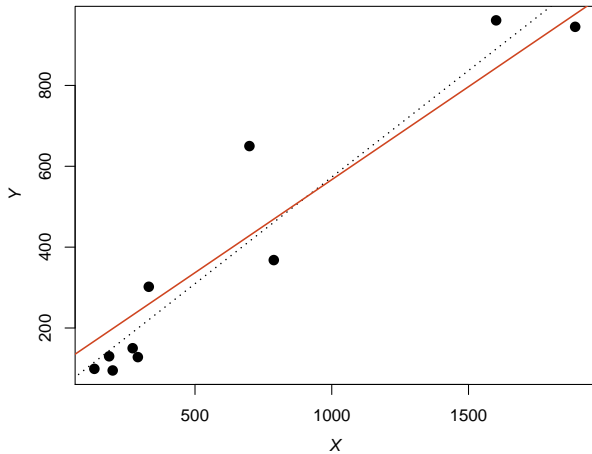
Szacowanie czasu programistów

iteracja 35



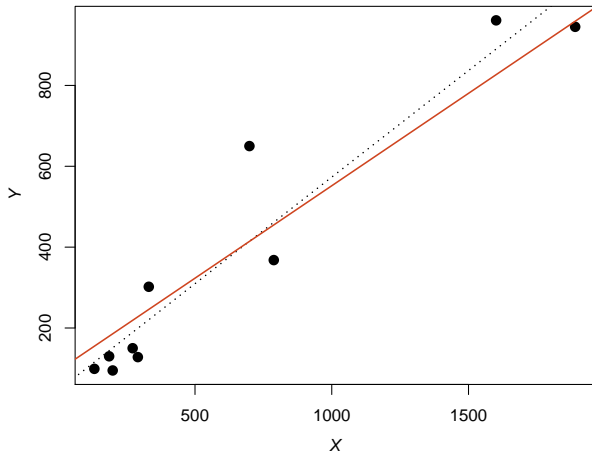
Szacowanie czasu programistów

iteracja 40



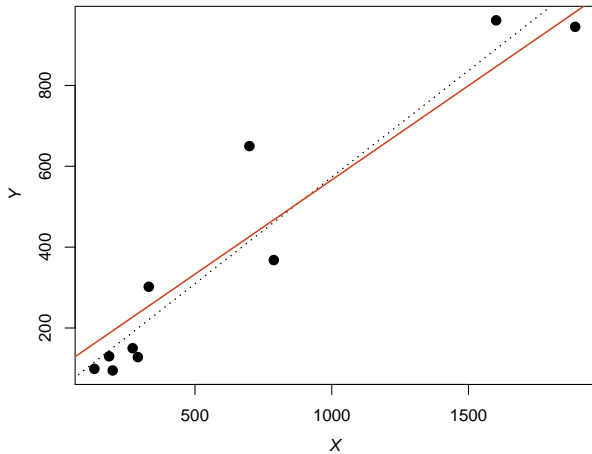
Szacowanie czasu programistów

iteracja 45



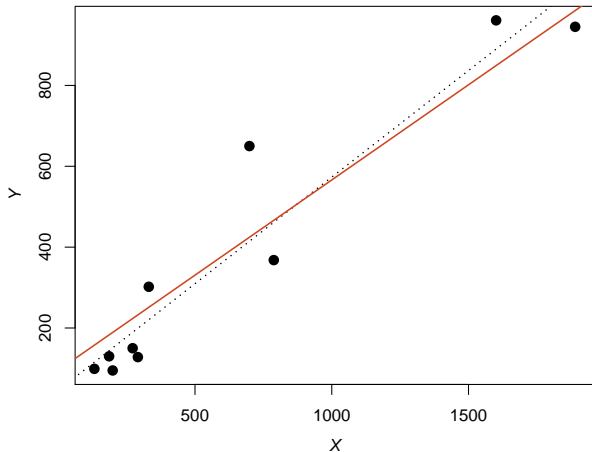
Szacowanie czasu programistów

iteracja 50



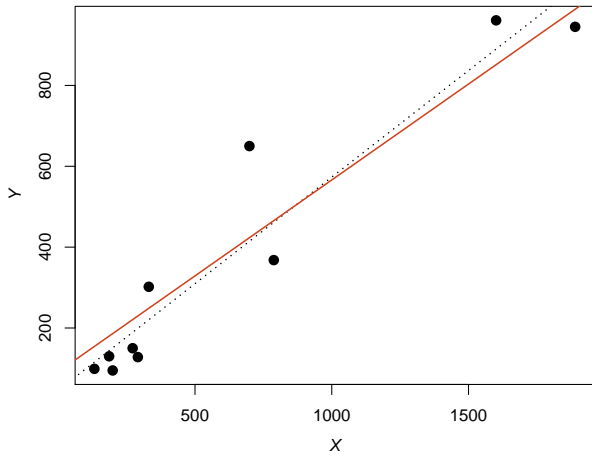
Szacowanie czasu programistów

iteracja 60



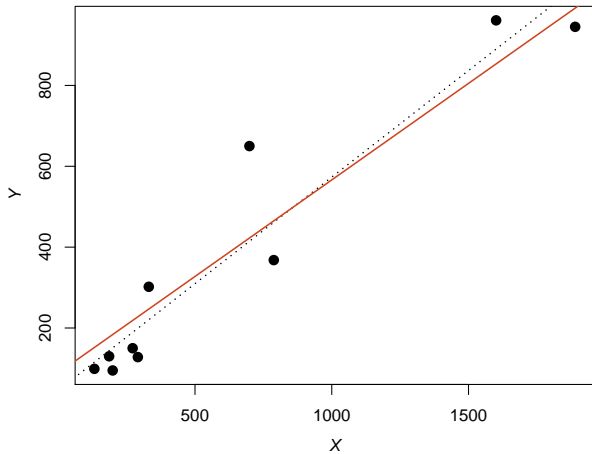
Szacowanie czasu programistów

iteracja 70



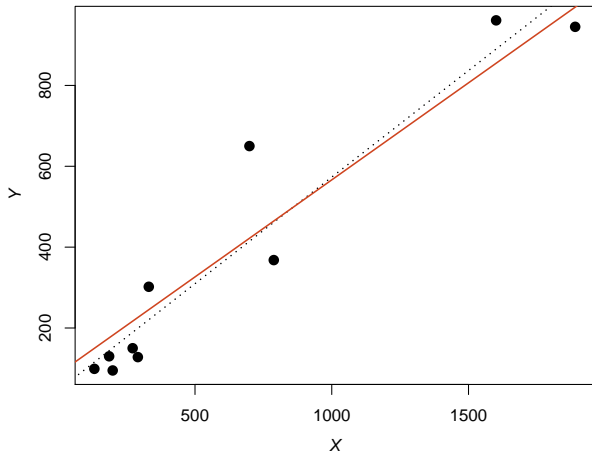
Szacowanie czasu programistów

iteracja 80



Szacowanie czasu programistów

iteracja 90



Szacowanie czasu programistów

iteracja 100

