

# Techniki Optymalizacji: Stochastyczny spadek wzdłuż gradientu I

Wojciech Kotłowski

Instytut Informatyki Politechniki Poznańskiej  
email: imię.nazwisko@cs.put.poznan.pl

pok. 2 (CW) tel. (61)665-2936 konsultacje: piątek 15:10-16:40  
Slajdy dostępne pod adresem: <http://www.cs.put.poznan.pl/wkottowski/to/>

27.11.2018

- 1 Metoda stochastycznego spadku wzdłuż gradientu
- 2 Stochastyczny gradient dla regresji liniowej
- 3 Przykład zastosowania SGD

- 1 Metoda stochastycznego spadku wzdłuż gradientu
- 2 Stochastyczny gradient dla regresji liniowej
- 3 Przykład zastosowania SGD

# Metoda spadku wzdłuż gradientu (Cauchy'ego)

Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wyznaczamy gradient w punkcie  $\mathbf{w}_{k-1}$ ,  $\nabla_L(\mathbf{w}_{k-1})$ .
  - Robimy krok wzdłuż ujemnego gradientu:

$$\mathbf{w}_k = \mathbf{w}_{k-1} - \alpha_k \nabla_L(\mathbf{w}_{k-1}),$$

gdzie  $\alpha_k$  jest długością kroku ustaloną np. przez przeszukiwanie liniowe.

W większości problemów uczenia maszynowego funkcja celu ma następującą postać:

$$L(\mathbf{w}) = \sum_{i=1}^n \ell_i(\mathbf{w})$$

W większości problemów uczenia maszynowego funkcja celu ma następującą postać:

$$L(\mathbf{w}) = \sum_{i=1}^n \ell_i(\mathbf{w})$$

- $L$  jest sumarycznym błędem na zbiorze uczącym.

W większości problemów uczenia maszynowego funkcja celu ma następującą postać:

$$L(\mathbf{w}) = \sum_{i=1}^n \ell_i(\mathbf{w})$$

- $L$  jest sumarycznym błędem na zbiorze uczącym.
- $\ell_i$  to błędy na poszczególnych obserwacjach.

# Specyficzna postać funkcji celu – przykłady



- Regresja liniowa – metoda najmniejszych kwadratów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

# Specyficzna postać funkcji celu – przykłady

- Regresja liniowa – metoda najmniejszych kwadratów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Regresja liniowa – min. wartości bezwzględnych błędów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{|y_i - \mathbf{w}^\top \mathbf{x}_i|}_{\ell_i(\mathbf{w})}$$

# Specyficzna postać funkcji celu – przykłady

- Regresja liniowa – metoda najmniejszych kwadratów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Regresja liniowa – min. wartości bezwzględnych błędów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{|y_i - \mathbf{w}^\top \mathbf{x}_i|}_{\ell_i(\mathbf{w})}$$

- Klasyfikacja liniowa – regresja logistyczna:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{\log \left( 1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)}_{\ell_i(\mathbf{w})}$$

# Specyficzna postać funkcji celu – przykłady

- Regresja liniowa – metoda najmniejszych kwadratów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Regresja liniowa – min. wartości bezwzględnych błędów

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{|y_i - \mathbf{w}^\top \mathbf{x}_i|}_{\ell_i(\mathbf{w})}$$

- Klasyfikacja liniowa – regresja logistyczna:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{\log \left( 1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i) \right)}_{\ell_i(\mathbf{w})}$$

- Klasyfikacja liniowa – funkcja zawiasowa:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{\left( 1 - y_i \mathbf{w}^\top \mathbf{x}_i \right)_+}_{\ell_i(\mathbf{w})}$$

# Metoda stochastycznego spadku wzdłuż gradientu

## Idea

Zamiast obliczać gradient na całej funkcji  $L$ , w danym kroku oblicz gradient tylko na pojedynczym elemencie  $\ell_i$ .

## Idea

Zamiast obliczać gradient na całej funkcji  $L$ , w danym kroku oblicz gradient tylko na pojedynczym elemencie  $l_i$ .

- Skąd nazwa „stochastyczny”? Ponieważ oryginalnie wybiera się element  $l_i$  losowo. . .

## Idea

Zamiast obliczać gradient na całej funkcji  $L$ , w danym kroku oblicz gradient tylko na pojedynczym elemencie  $l_i$ .

- Skąd nazwa „stochastyczny”? Ponieważ oryginalnie wybiera się element  $l_i$  losowo. . .
- . . . ale w praktyce zwykle przechodzi się po całym zbiorze danych w losowej kolejności.

Minimalizacja funkcji  $L(\mathbf{w})$ :



Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .

Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności

Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .

Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .
  - Wyznaczamy gradient funkcji  $\ell_i$  w punkcie  $\mathbf{w}_{k-1}$ ,  $\nabla_{\ell_i}(\mathbf{w}_{k-1})$ .

# Metoda stochastycznego spadku wzdłuż gradientu

Minimalizacja funkcji  $L(\mathbf{w})$ :

- 1 Zaczynamy od wybranego rozwiązania startowego, np.  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .
  - Wyznaczamy gradient funkcji  $\ell_i$  w punkcie  $\mathbf{w}_{k-1}$ ,  $\nabla_{\ell_i}(\mathbf{w}_{k-1})$ .
  - Robimy krok wzdłuż negatywnego gradientu:

$$\mathbf{w}_k = \mathbf{w}_{k-1} - \alpha_k \nabla_{\ell_i}(\mathbf{w}_{k-1}),$$

gdzie  $\alpha_k$  jest długością kroku.

# Zalety metody stochastycznego spadku wzdłuż gradientu

# Zalety metody stochastycznego spadku wzdłuż gradientu

- **Szybkość:** obliczenie gradientu wymaga wzięcia tylko jednej obserwacji.

# Zalety metody stochastycznego spadku wzdłuż gradientu

- **Szybkość**: obliczenie gradientu wymaga wzięcia tylko jednej obserwacji.
- **Skalowalność**: cały zbiór danych nie musi nawet znajdować się w pamięci operacyjnej.



# Zalety metody stochastycznego spadku wzdłuż gradientu

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

# Zalety metody stochastycznego spadku wzdłuż gradientu

- **Szybkość**: obliczenie gradientu wymaga wzięcia tylko jednej obserwacji.
- **Skalowalność**: cały zbiór danych nie musi nawet znajdować się w pamięci operacyjnej.
- **Prostota**: gradient funkcji  $\ell_i$  daje bardzo prosty wzór na modyfikację wag.
- **Stochastyczny gradient jest obecnie najbardziej popularną metodą stosowaną w uczeniu maszynowym!**  
(*stochastic gradient descent* – SGD).

# Wady metody stochastycznego spadku wzdłuż gradientu

- **Wolna zbieżność**: czasem gradient stochastyczny zbiega wolno i wymaga wielu iteracji po zbiorze uczącym.

- **Wolna zbieżność**: czasem gradient stochastyczny zbiega wolno i wymaga wielu iteracji po zbiorze uczącym.
- **Problem z ustaleniem długości kroku  $\alpha_k$** : wyznaczenie  $\alpha_k$  przez przeszukiwanie liniowe nie przynosi dobrych rezultatów, ponieważ optymalizujemy oryginalnej funkcji  $L$  tylko jej jeden składnik  $\ell_i$ .

- **Wolna zbieżność**: czasem gradient stochastyczny zbiega wolno i wymaga wielu iteracji po zbiorze uczącym.
- **Problem z ustaleniem długości kroku  $\alpha_k$** : wyznaczenie  $\alpha_k$  przez przeszukiwanie liniowe nie przynosi dobrych rezultatów, ponieważ optymalizujemy oryginalnej funkcji  $L$  tylko jej jeden składnik  $\ell_i$ .
- **Zalety znacznie przewyższają wady!**



- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności.



- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności.
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).

- Zwykle nie losuje się obserwacji, ale przechodzi się po zbiorze danych w losowej kolejności.
- Zbieżność wymaga często przejścia parokrotnie po całym zbiorze danych (jednokrotne przejście nazywa się **epoką**).
- Metody ustalania współczynników długości kroku  $\alpha_k$ :

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

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

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

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

- 1 Metoda stochastycznego spadku wzdłuż gradientu
- 2 Stochastyczny gradient dla regresji liniowej
- 3 Przykład zastosowania SGD

# Regresja liniowa – metoda najmniejszych kwadratów (LS)



- Funkcja celu:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Funkcja celu:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Pochodne funkcji  $\ell_i(\mathbf{w})$ :

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

- Funkcja celu:

$$L(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}_{\ell_i(\mathbf{w})}$$

- Pochodne funkcji  $\ell_i(\mathbf{w})$ :

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

- Gradient funkcji  $\ell_i(\mathbf{w})$ :

$$\nabla_{\ell_i}(\mathbf{w}) = -2(y_i - \mathbf{w}^\top \mathbf{x}_i)\mathbf{x}_i$$

# SGD dla regresji liniowej LS

- 1 Zaczynamy od  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .
  - Wyznaczamy gradient  $\ell_i$  w punkcie  $\mathbf{w}_{k-1}$ ,  
 $-2(y_i - \mathbf{w}_{k-1}^\top \mathbf{x}_i) \mathbf{x}_i$ .
  - Robimy krok wzdłuż negatywnego gradientu:

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

# SGD dla regresji liniowej LS

- 1 Zaczynamy od  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .
  - Wyznaczamy gradient  $\ell_i$  w punkcie  $\mathbf{w}_{k-1}$ ,  
 $-2(y_i - \mathbf{w}_{k-1}^\top \mathbf{x}_i) \mathbf{x}_i$ .
  - Robimy krok wzdłuż negatywnego gradientu:

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

- 
- Krok w kierunku wektora  $\mathbf{x}_i$ , z długością kroku równą  $2\alpha_k(y_i - \mathbf{w}_{k-1}^\top \mathbf{x}_i)$ .

# SGD dla regresji liniowej LS

- 1 Zaczynamy od  $\mathbf{w}_0 = \mathbf{0}$ .
- 2 Dla  $k = 1, 2, \dots$  aż do zbieżności
  - Wylosuj  $i \in \{1, \dots, n\}$ .
  - Wyznaczamy gradient  $\ell_i$  w punkcie  $\mathbf{w}_{k-1}$ ,  
 $-2(y_i - \mathbf{w}_{k-1}^\top \mathbf{x}_i) \mathbf{x}_i$ .
  - Robimy krok wzdłuż negatywnego gradientu:

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

- 
- Krok w kierunku wektora  $\mathbf{x}_i$ , z długością kroku równą  $2\alpha_k(y_i - \mathbf{w}_{k-1}^\top \mathbf{x}_i)$ .
  - Długość kroku proporcjonalna do „przeszacowania predykcji”.

- 1 Metoda stochastycznego spadku wzdłuż gradientu
- 2 Stochastyczny gradient dla regresji liniowej
- 3 Przykład zastosowania SGD

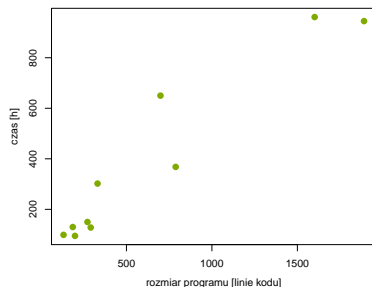
## Przykład: szacowanie czasu pracy programistów

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



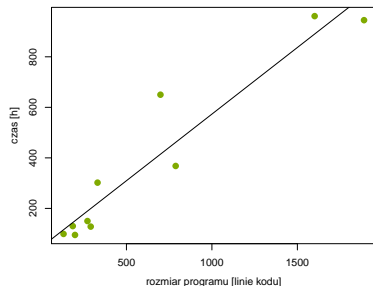
# Przykład: szacowanie czasu pracy programistów

| $X$              | $Y$             |
|------------------|-----------------|
| 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             |

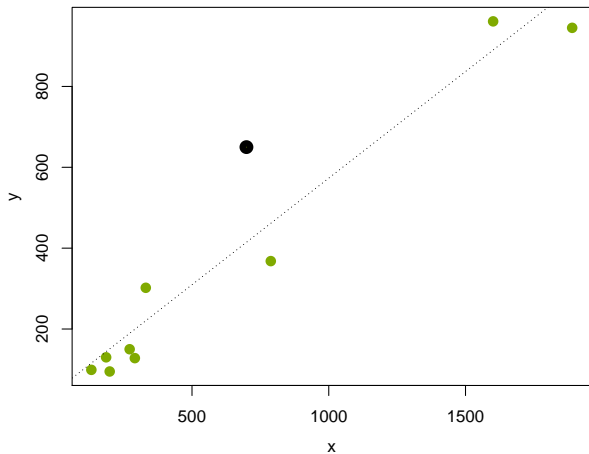


$$w_0 = 45.93, w_1 = 0.5273$$

- 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 istotne dla zbieżności metody!)

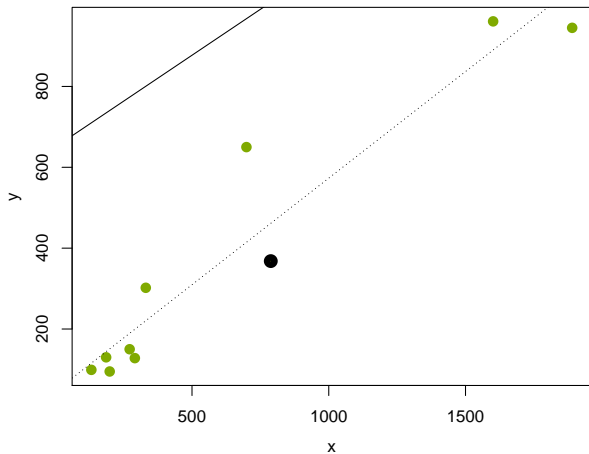
# Szacowanie Czasu programistów

iteracja 1



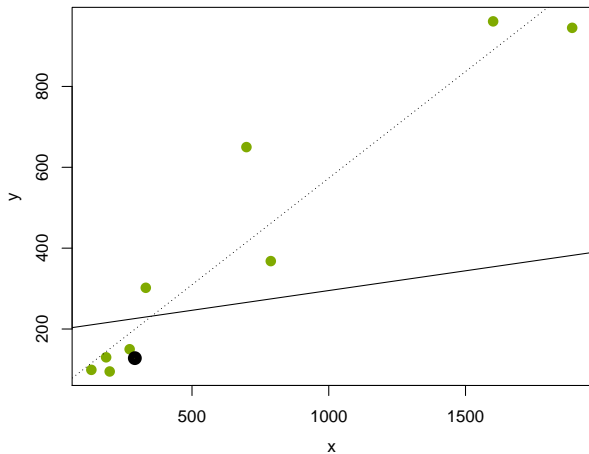
# Szacowanie Czasu programistów

iteracja 2



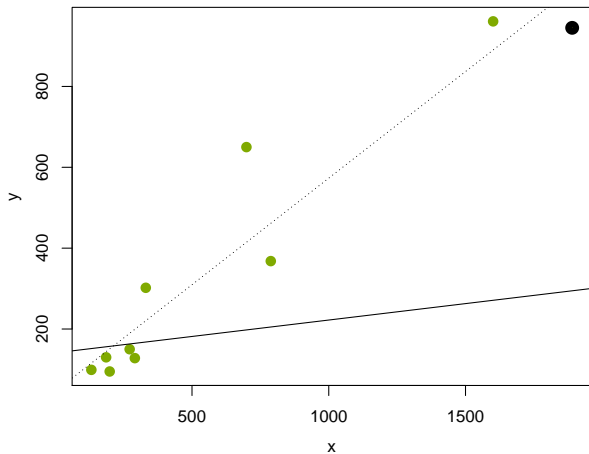
# Szacowanie Czasu programistów

iteracja 3

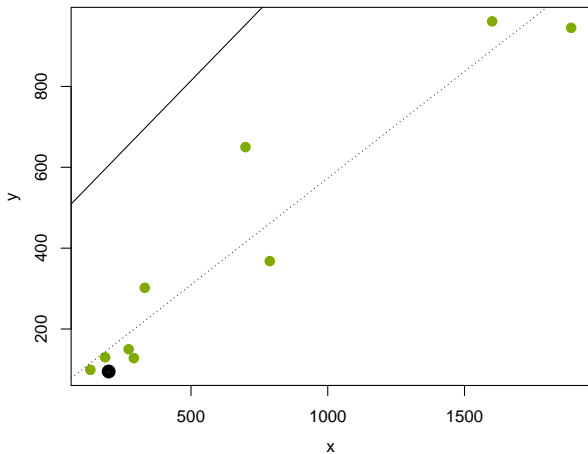


# Szacowanie Czasu programistów

iteracja 4



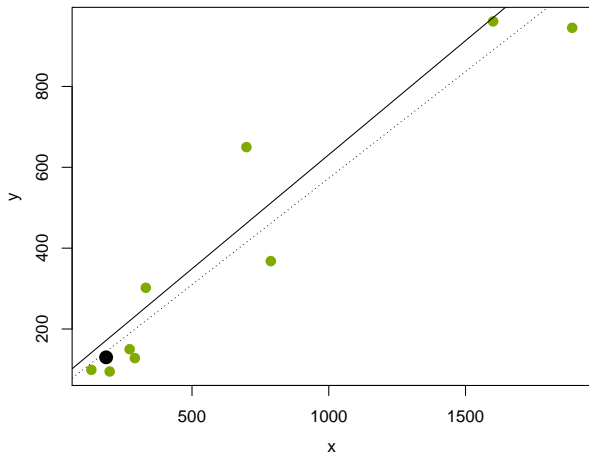
iteracja 5



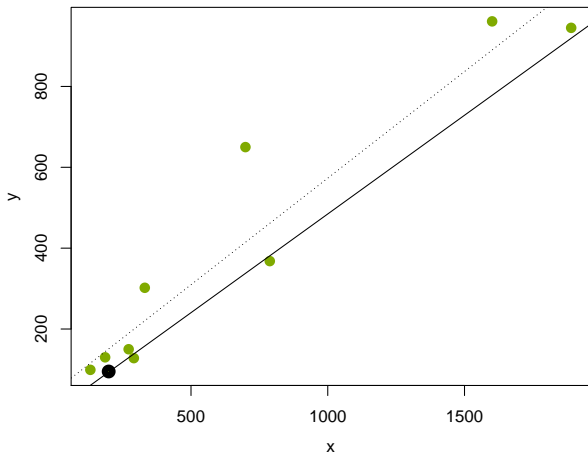


# Szacowanie Czasu programistów

iteracja 10

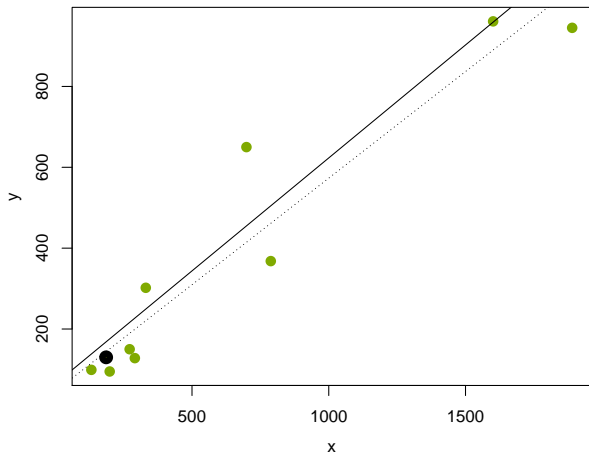


iteracja 15



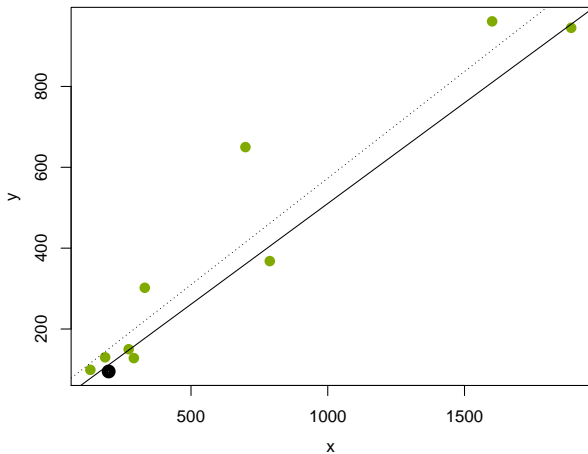
# Szacowanie Czasu programistów

iteracja 20



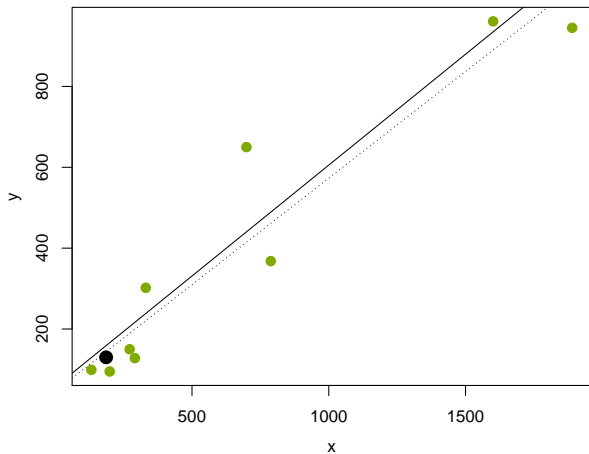
# Szacowanie Czasu programistów

iteracja 25

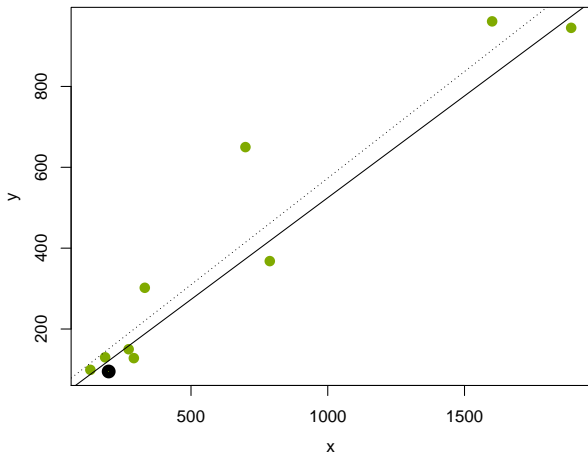


# Szacowanie Czasu programistów

iteracja 30

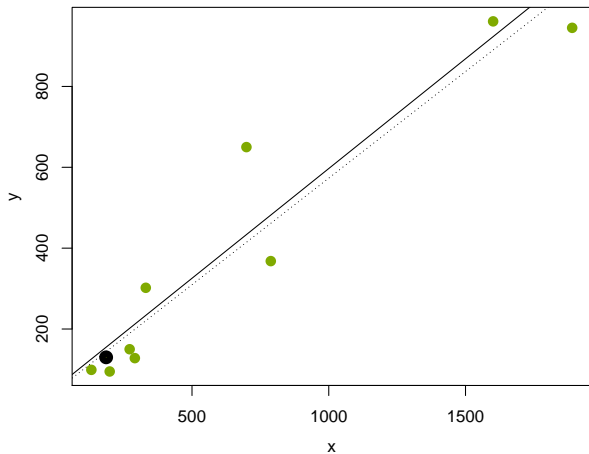


iteracja 35



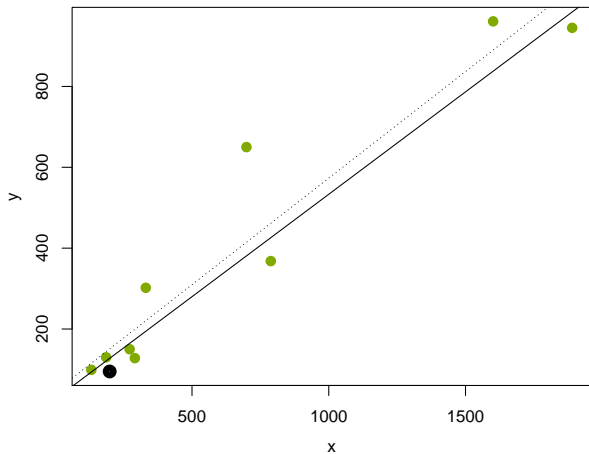
# Szacowanie Czasu programistów

iteracja 40



# Szacowanie Czasu programistów

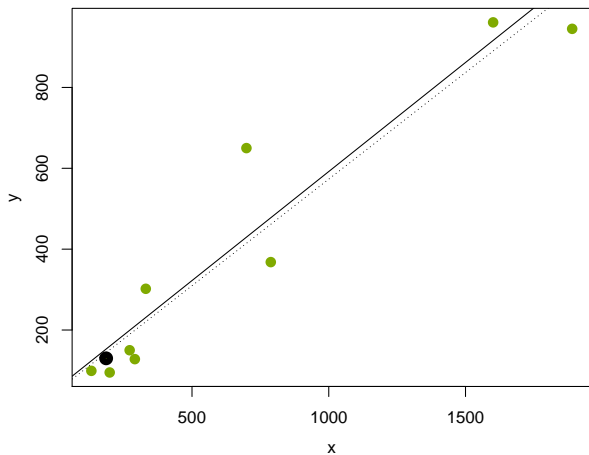
iteracja 45





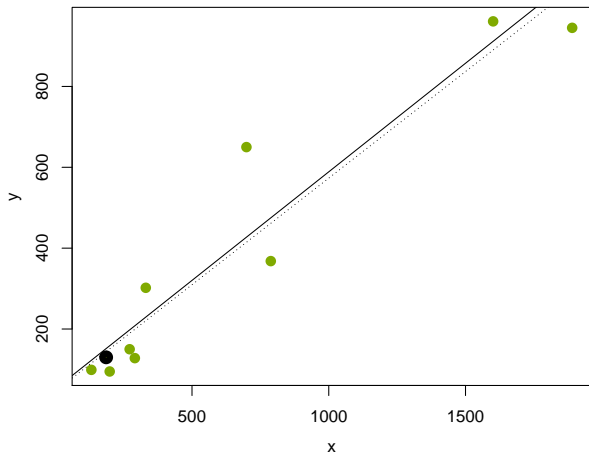
# Szacowanie Czasu programistów

iteracja 50

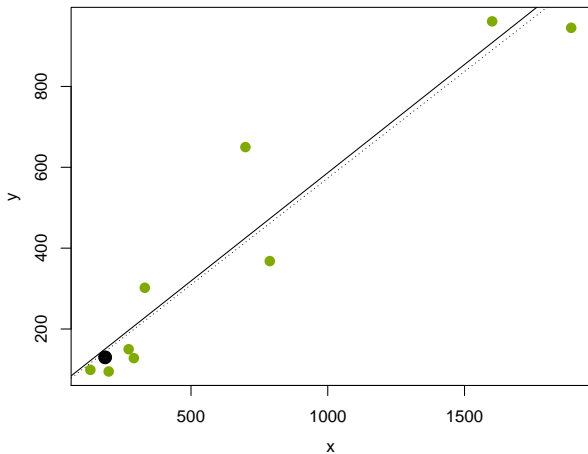


# Szacowanie Czasu programistów

iteracja 60

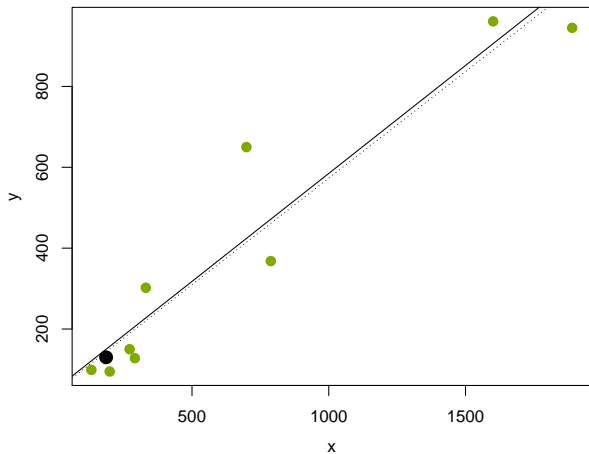


iteracja 70



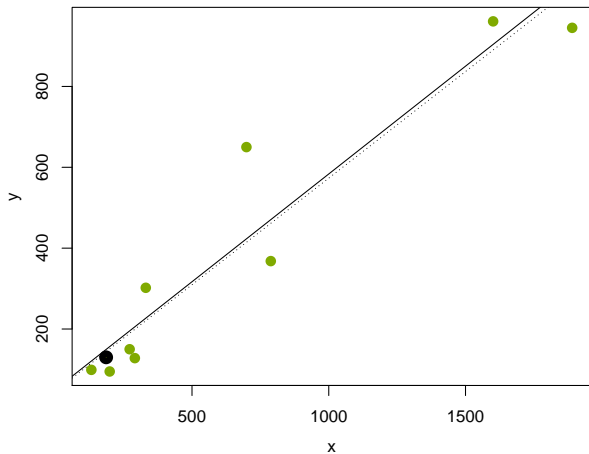
# Szacowanie Czasu programistów

iteracja 80



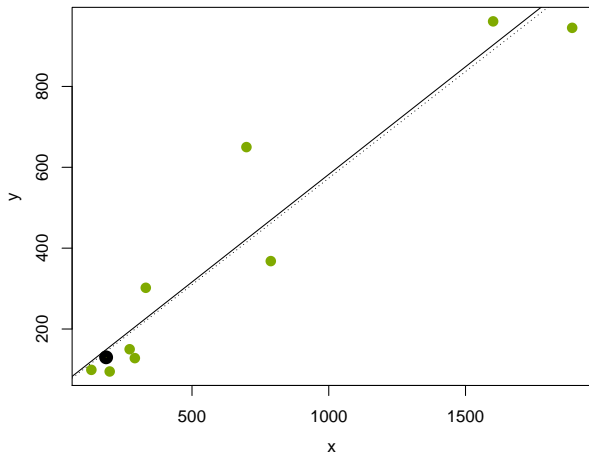
# Szacowanie Czasu programistów

iteracja 90



# Szacowanie Czasu programistów

iteracja 100



Koniec na dzisiaj :)