

# Regresja logistyczna, problemy bias/variance, diagnostyka



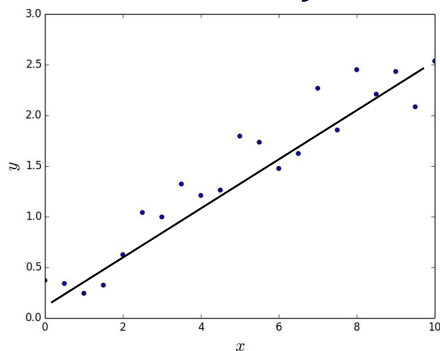
# 1. Regresja logistyczna

# Zadania dla uczenia maszynowego

Punkt wyjścia:

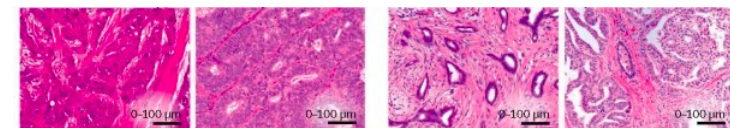
- dane wejściowe (treningowe):  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
- oczekiwane wyniki:  $\mathbf{y}$
- liczba próbek danych treningowych:  $m$

## REGRESJA



## KLASYFIKACJA

- rozpoznawanie danych: obiektów graficznych, dźwiękowych, danych medycznych, personalnych, identyfikacja spamu, zagrożeń sieciowych



1. hipoteza:

$$h(\mathbf{x}) = \theta_0 \cdot \mathbf{x}_0 + \theta_1 \cdot \mathbf{x}_1 + \dots + \theta_n \cdot \mathbf{x}_n$$

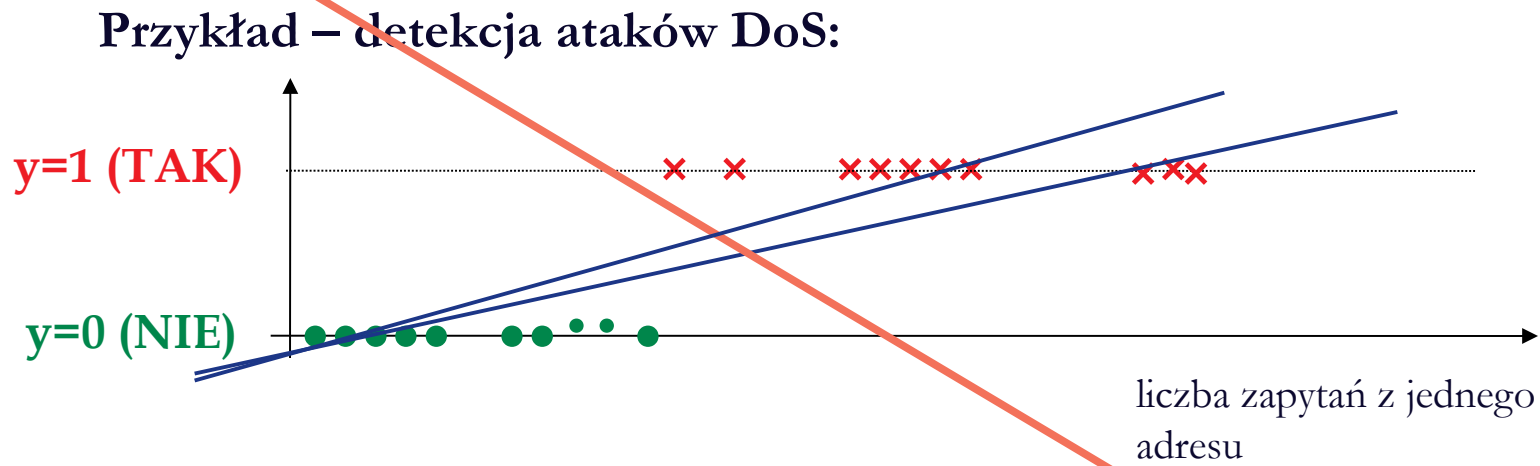
2. funkcja kosztu:

$$J = \frac{1}{2m} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

3. gradient descent:

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)}) \cdot \mathbf{x}_j^{(i)}$$

# Regresja liniowa?



Hipoteza:  $h(\mathbf{x}) = \theta_0 \cdot \mathbf{x}_0 + \theta_1 \cdot \mathbf{x}_1 + \dots + \theta_n \cdot \mathbf{x}_n ?$

# Regresja logistyczna

Hipoteza w regresji liniowej:

$$h(\mathbf{x}) = \theta_0 \cdot \mathbf{x}_0 + \theta_1 \cdot \mathbf{x}_1 + \dots + \theta_n \cdot \mathbf{x}_n = \boldsymbol{\theta} \mathbf{x}^T$$

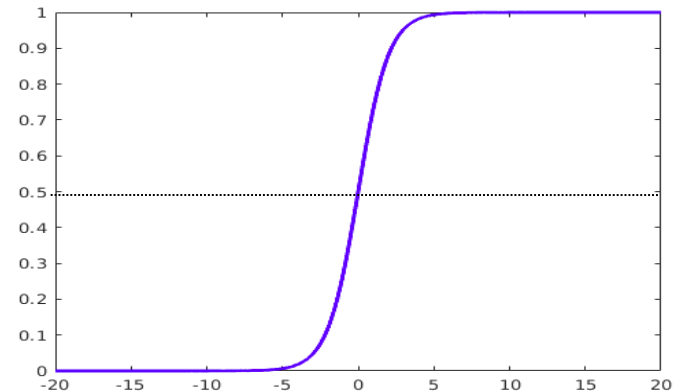
$$\mathbf{x} = (\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_n)$$

Hipoteza w regresji logistycznej:

$$h(\mathbf{x}) = g(\boldsymbol{\theta} \mathbf{x}^T) = \frac{1}{1 + \exp(-\boldsymbol{\theta} \mathbf{x}^T)}$$

funkcja sigmoidalna (logistyczna):



Interpretacja:

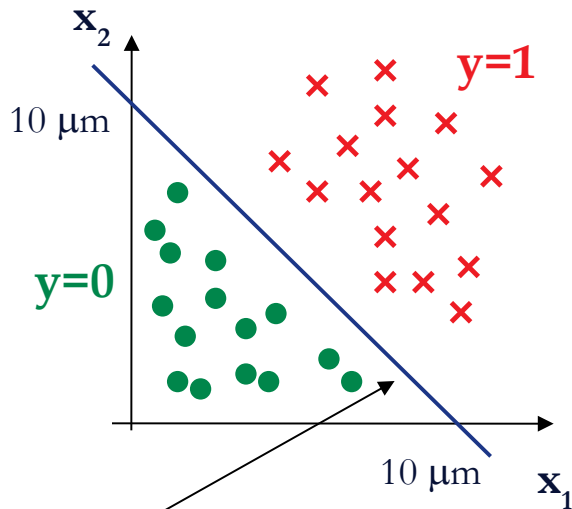
$h(\mathbf{x}) \rightarrow$  prawdopodobieństwo, że  $y(\mathbf{x}) = 1$

$\boldsymbol{\theta} \mathbf{x}^T > 0 \rightarrow p > 50\%$ , że  $y(\mathbf{x}) = 1$

$$g(z) = 1 / (1 + \exp(-z))$$

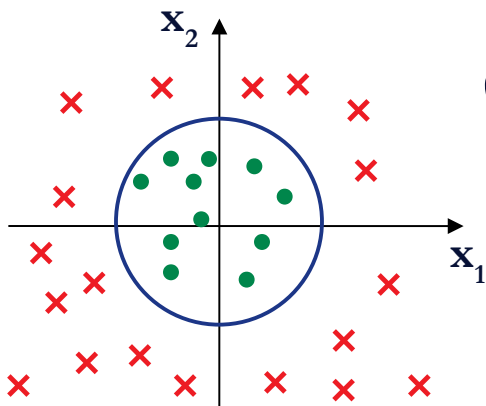
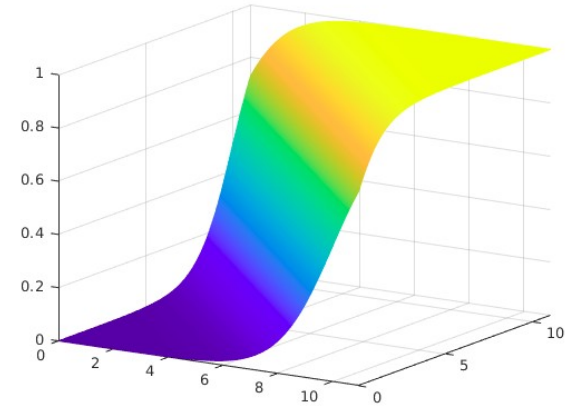
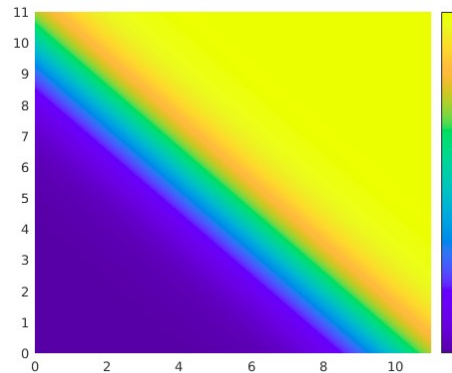
# Przykłady regresji logistycznej

Przykład – detekcja bakterii:



$$\theta \mathbf{x}^T = x_1 + x_2 - 10$$

$$h(\mathbf{x}) = \frac{1}{1 + \exp(-\theta \mathbf{x}^T)}$$

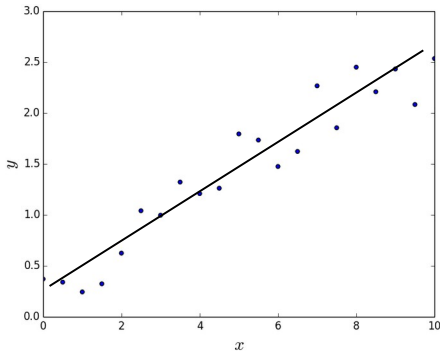


$$\theta \mathbf{x}^T = x_1^2 + x_2^2 - r^2$$

$$h(\mathbf{x}) = \frac{1}{1 + \exp(-x_1^2 - x_2^2 + r^2)}$$

# Funkcja kosztu

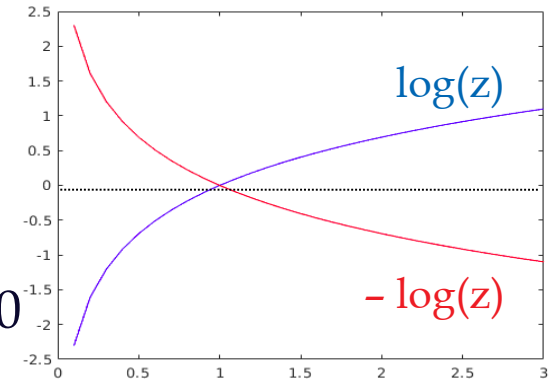
Regresja liniowa:



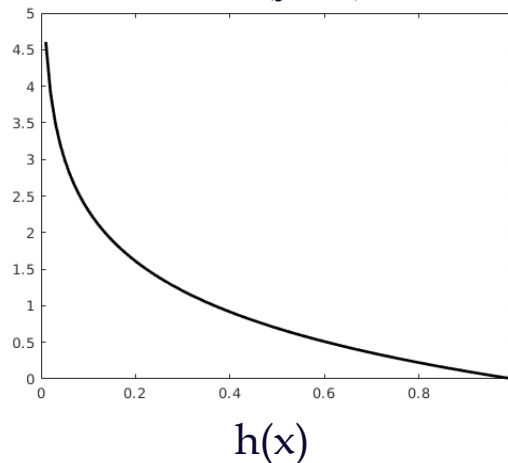
$$J = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

Propozycja  
dla regresji logistycznej:

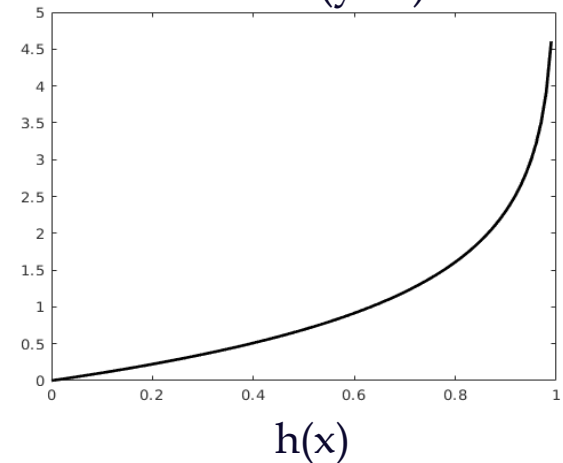
$$\text{Koszt} = \begin{cases} -\log(h(x)), & \text{dla } y=1 \\ -\log(1-h(x)), & \text{dla } y=0 \end{cases}$$



Koszt (y=1)



Koszt (y=0)



$$J = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \cdot \log(h(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-h(x^{(i)}))$$

# Gradient descent dla r. logistycznej

$$\text{Koszt} = \begin{cases} -\log(h(x)), & \text{dla } y=1 \\ -\log(1-h(x)), & \text{dla } y=0 \end{cases}$$
$$h(x) = \frac{1}{1 + \exp(-\theta x^T)}$$

$$J = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \cdot \log(h(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-h(x^{(i)}))$$

0 lub 1

Cel algorytmu gradient descent: policzyć  $\min_{\theta_0, \theta_1, \dots, \theta_n} J(\theta)$

Dla każdej  $\theta_j, j=0,1,\dots,n$ , powtarzamy:

$$\theta_j = \theta_j - \alpha \frac{d}{d\theta_j} J(\theta)$$

... (liczymy pochodną)

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

Wygląda znajomo :-)

# Jak liczyć błąd dla r. logistycznej?

$$h(\mathbf{x}) = \frac{1}{1 + \exp(-\theta \mathbf{x}^T)}$$
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \cdot \log(h(\mathbf{x}^{(i)})) + (1-y^{(i)}) \cdot \log(1-h(\mathbf{x}^{(i)}))$$

0 lub 1

Chcemy policzyć jak wiele klasyfikacji się nie udaje, a więc:

$$\text{err}(h(\mathbf{x}), y) = \begin{cases} 0, & \text{gdy } h(\mathbf{x}) \geq 0.5 \text{ i } y=1 \text{ oraz gdy } h(\mathbf{x}) < 0.5 \text{ i } y=0 \\ 1, & \text{w przeciwnym wypadku} \end{cases}$$

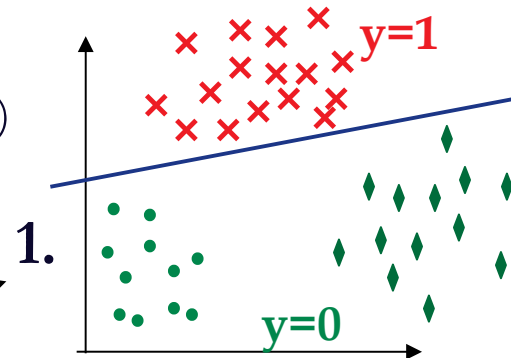
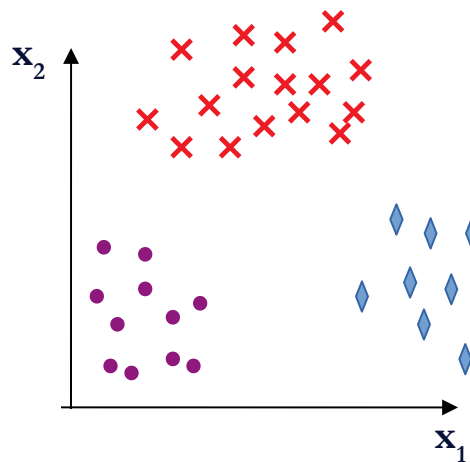
$$\text{Test error} = \frac{1}{m} \sum_{i=1}^m \text{err}(h(\mathbf{x}^{(i)}), y^{(i)})$$

# A jeżeli mamy więcej kategorii?

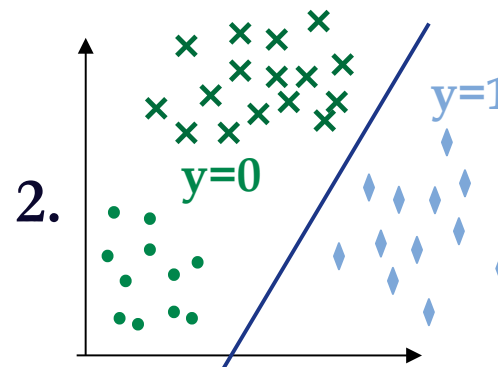
*multi-class classification*

## Algorytm one-vs-all:

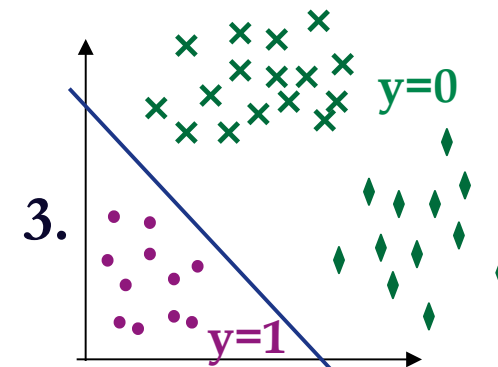
(równoległe trenujemy kilka klasyfikatorów)



→  $h_1(\mathbf{x}) =$   
p-stwo, że obiekt  
to  $\times$



→  $h_2(\mathbf{x}) =$   
p-stwo, że obiekt  
to  $\diamond$



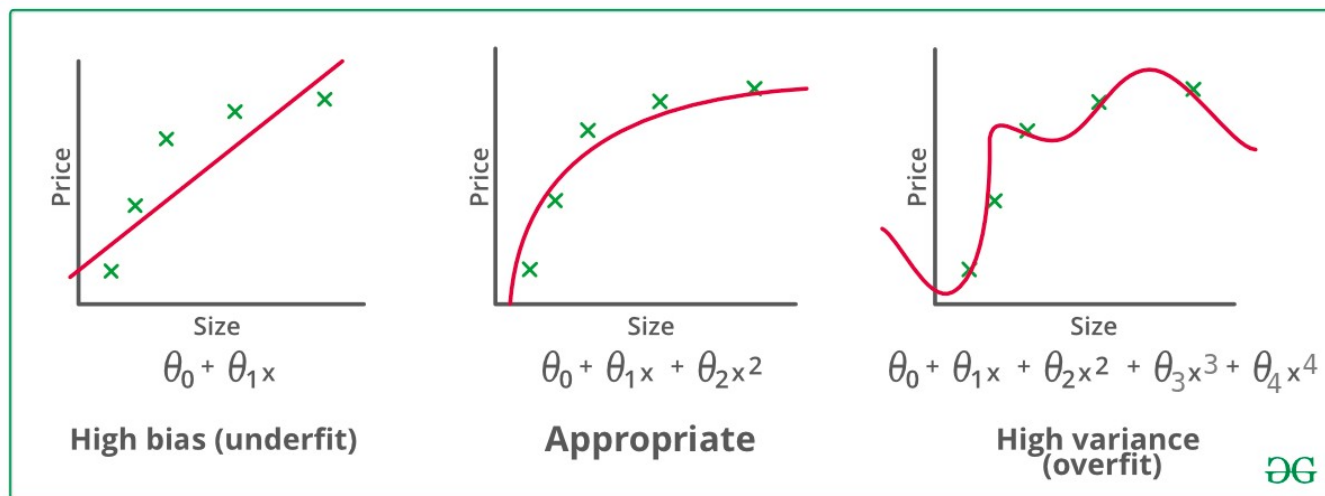
→  $h_3(\mathbf{x}) =$   
p-stwo, że obiekt  
to  $\bullet$

## **2. Problem of bias/variance**

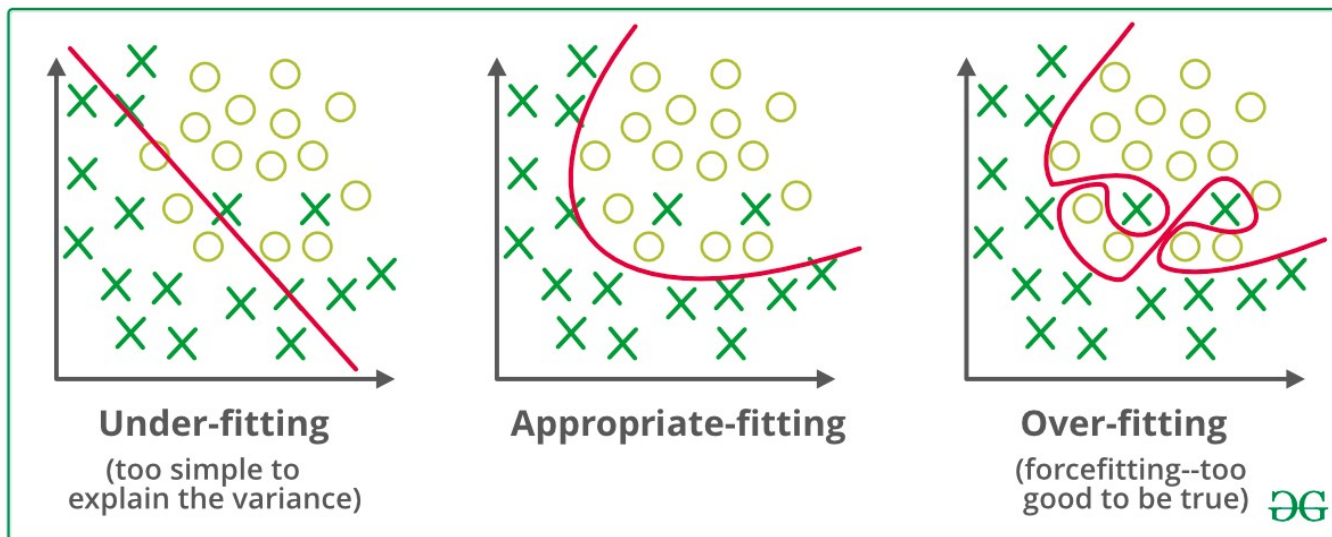
# High bias/variance

*underfitting / overfitting*

←  
**Regresja liniowa:**



**Regresja logistyczna:**



# Underfitting – co robić?

---

## 1. zwiększyć wymiar danych

(dodać nowe cechy – kolumny wśród danych wejściowych)

## 2. zwiększyć wymiar danych

(dodać wielomiany lub inne funkcje już istniejących danych)

**(model jest zbyt prosty)**

# Overfitting – co robić?

---

## 1. zmniejszyć wymiar danych

(usunąć wybrane cechy – kolumny wśród danych wejściowych)

## 2. zastosować **regularyzację**:

(zmniejszanie współczynników  $\theta$ )

**(mamy za mało próbek danych)**

# Regularyzacja dla r. liniowej

Hipoteza:

$$h(\mathbf{x}) = \theta_0 \cdot \mathbf{x}_0 + \theta_1 \cdot \mathbf{x}_1 + \dots + \theta_n \cdot \mathbf{x}_n$$

Funkcja kosztu:

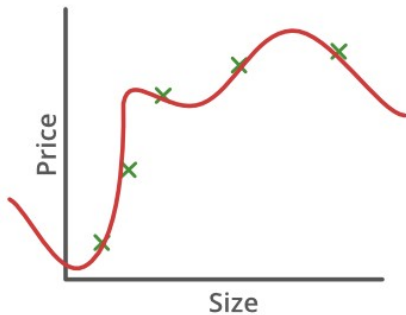
$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Gradient descent:

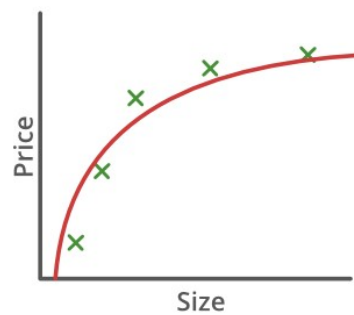
$$\theta_j = \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

trochę mniejsze niż 1

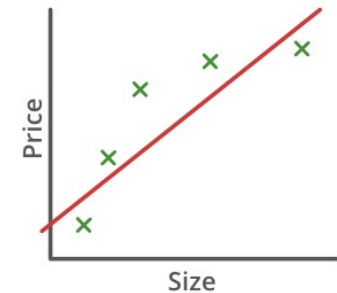
współczynnik regularyzacji



$\lambda$  zbyt mała lub 0



$\lambda$  "w sam raz"



zbyt duża  $\lambda$

# Regularyzacja dla r. logistycznej

Hipoteza:

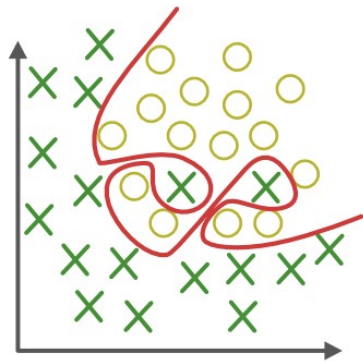
$$h(\mathbf{x}) = \frac{1}{1 + \exp(-\theta \mathbf{x}^T)} \quad \theta \mathbf{x}^T = \theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \dots + \theta_n \cdot x_n$$

Funkcja kosztu:

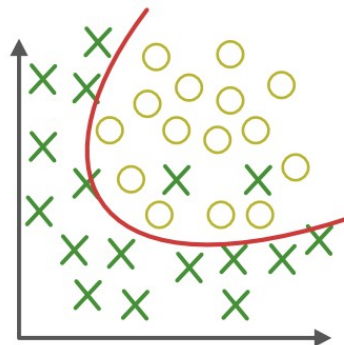
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \cdot \log(h(x^{(i)})) + (1-y^{(i)}) \cdot \log(1-h(x^{(i)}))] + \frac{1}{2m} \lambda \sum_{j=1}^n \theta_j^2$$

Gradient descent:

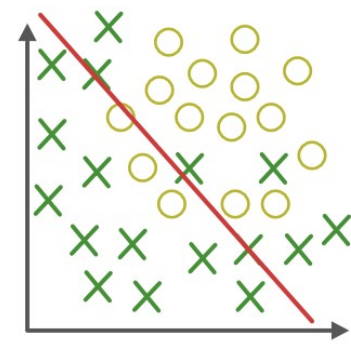
$$\theta_j = \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$



$\lambda$  zbyt mała lub 0



$\lambda$  "w sam raz"



zbyt duża  $\lambda$

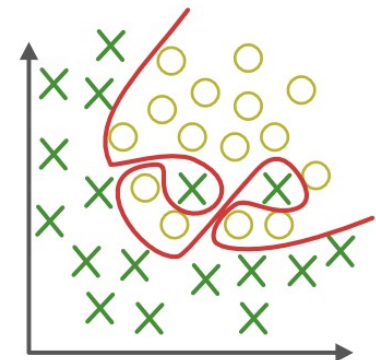
# 3. Diagnostyka

# Mam algorytm, ale nie działa dobrze...

← ... błędy są zbyt duże...

## Możliwe strategie:

- dodać więcej danych wejściowych (próbek)?
- zwiększyć wymiar danych (liczbę cech)?
- dodać nowe dane – wielomiany cech istniejących?
- zmniejszyć wymiar danych?
- regularyzacja – większa  $\lambda$ ?
- regularyzacja – mniejsza  $\lambda$ ?



**Wielowymiarowe dane  
nie dadzą się tak ładnie zwizualizować!**

# Diagnostyka – zbiory danych

Podział (losowy!) danych wejściowych:

$x_1$	$x_2$	$x_3$	$y$
6k	27	3.2	300
11.5k	39	4.3	500
8k	41	4.5	350
4k	84	7.8	650
15k	42	3.7	850
18.5k	36	4.1	800
22k	149	13.6	1800
21k	29	3.0	1000
23k	69	7.7	1400
24k	71	10.1	1450

$m$  próbek danych

zbiór treningowy ( $m_{\text{train}}$ ) ~60%  
*training set*

zbiór walidacyjny ( $m_{\text{cv}}$ ) ~20%  
*validation set*

zbiór testowy ( $m_{\text{test}}$ ) ~20%  
*test set*

# Po co nam 3 zbiory?

---

Uczenie modelu na zbiorze treningowym  $m_{\text{train}}$ :

- hipoteza  $h(x)$
- funkcja kosztu:  $J_{\text{train}}(\theta) = \frac{1}{2m_{\text{train}}} \left[ \sum_{i=1}^{m_{\text{train}}} (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$
- gradient descent: iteracyjna optymalizacja  $\theta$

Walidacja parametrów modelu na zbiorze walidacyjnym  $m_{\text{cv}}$ :

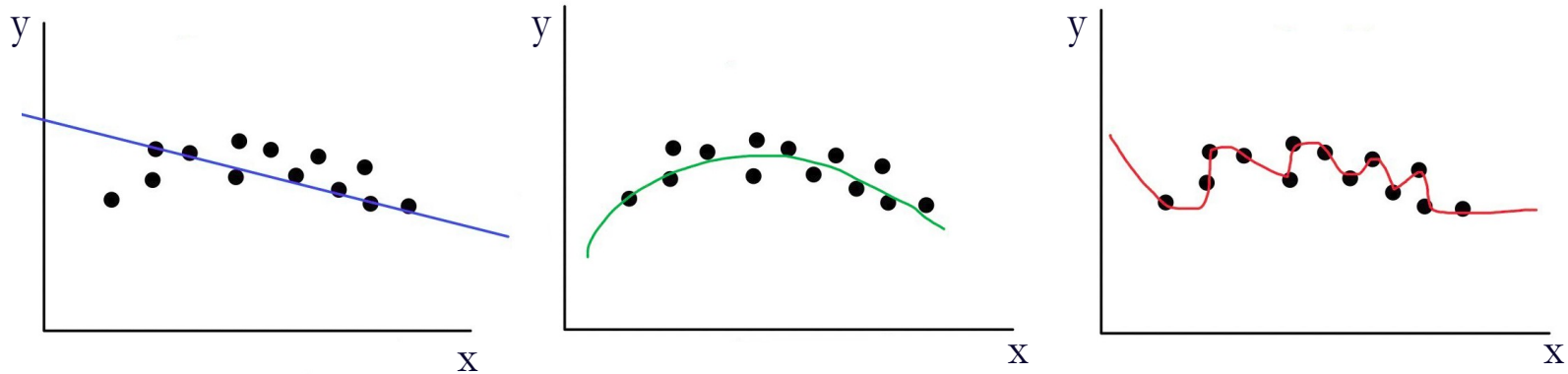
- funkcja kosztu:  $J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h(x^{(i)}) - y^{(i)})^2$

Ostateczny test modelu na niezależnym (\*) zbiorze testowym  $m_{\text{test}}$ :

- funkcja kosztu:  $J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h(x^{(i)}) - y^{(i)})^2$

# Wybór i walidacja parametrów modelu

←  
Przykład: regresja liniowa, dobór stopnia wielomianu:



$$\theta_0 + \theta_1 x$$

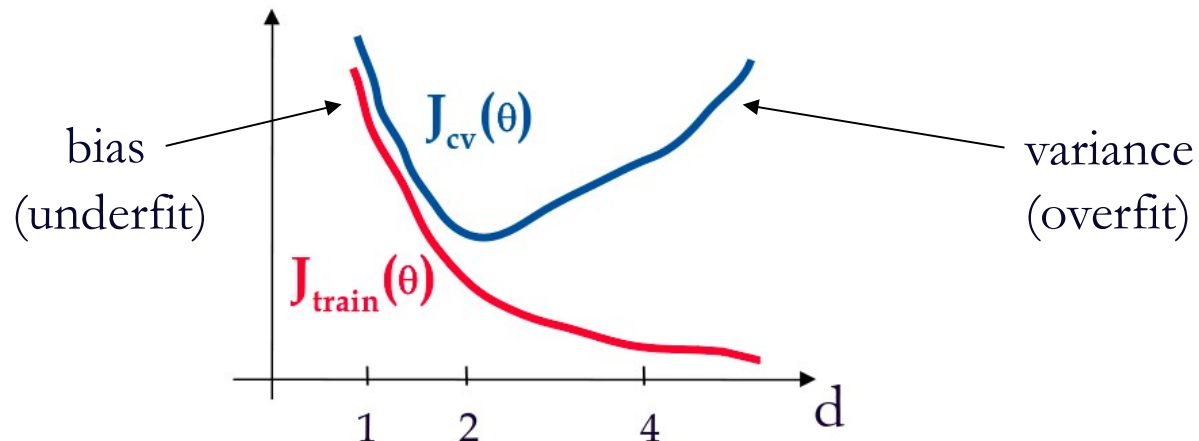
$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Najwyższy  
stopień  $\longrightarrow$   $d=1$   
wielomianu:

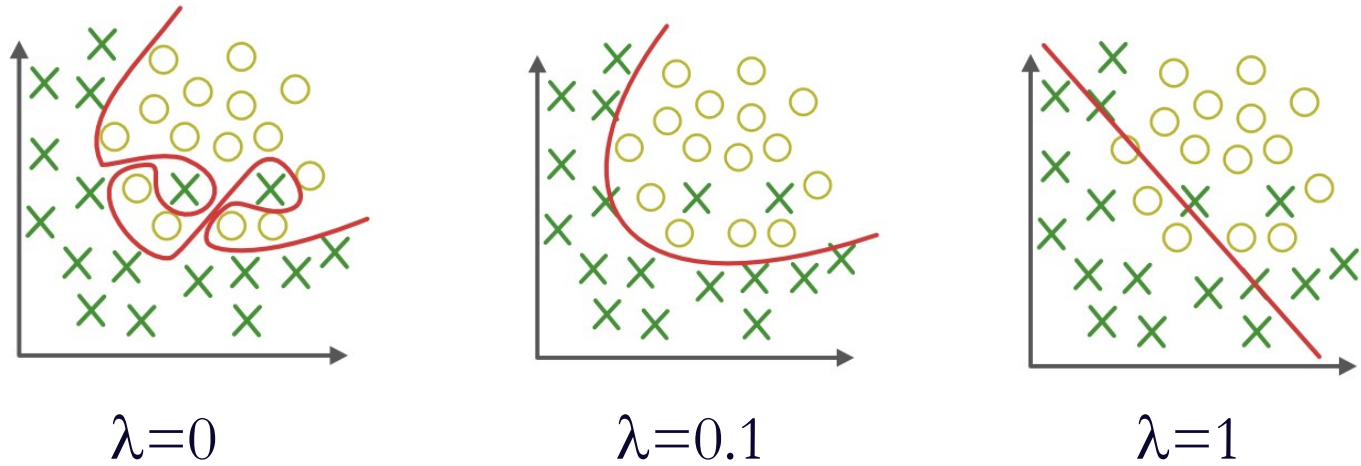
$d=2$

$d=4$

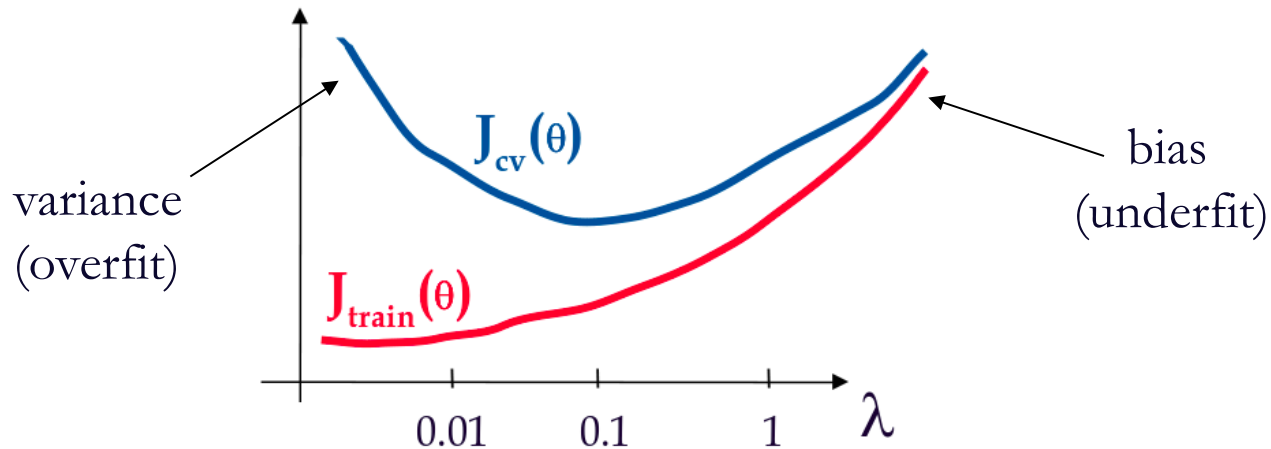


# Wybór i walidacja parametrów modelu

Przykład: regresja logistyczna z regularyzacją, dobór  $\lambda$  :



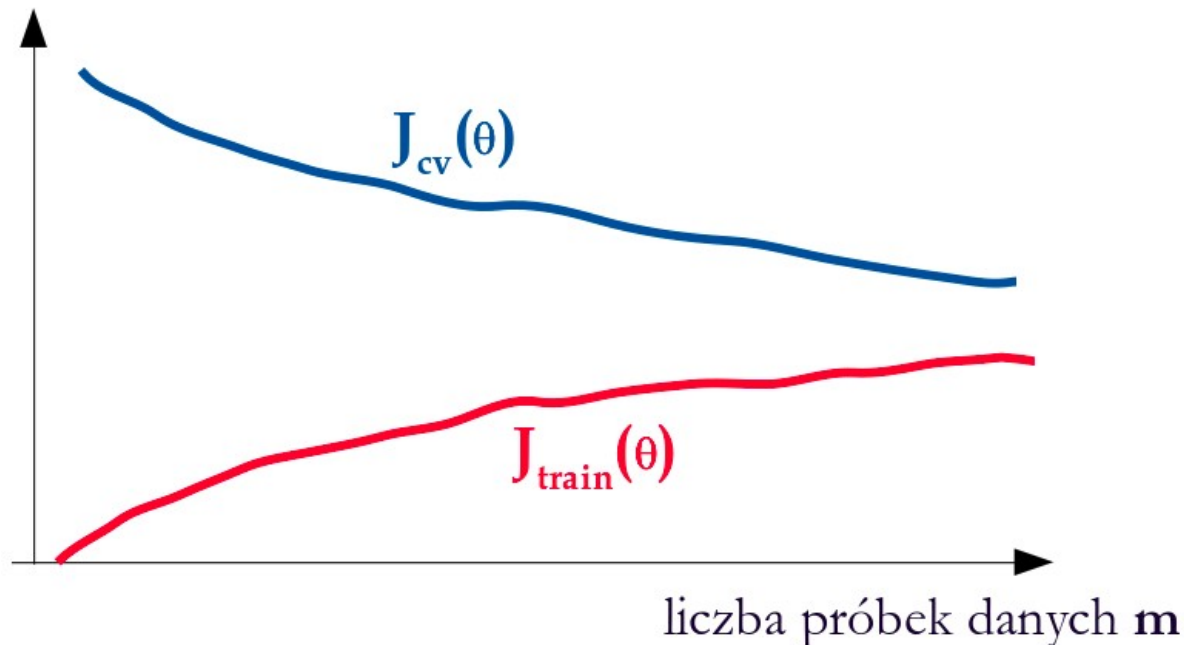
Parametr  
regularyzacji:



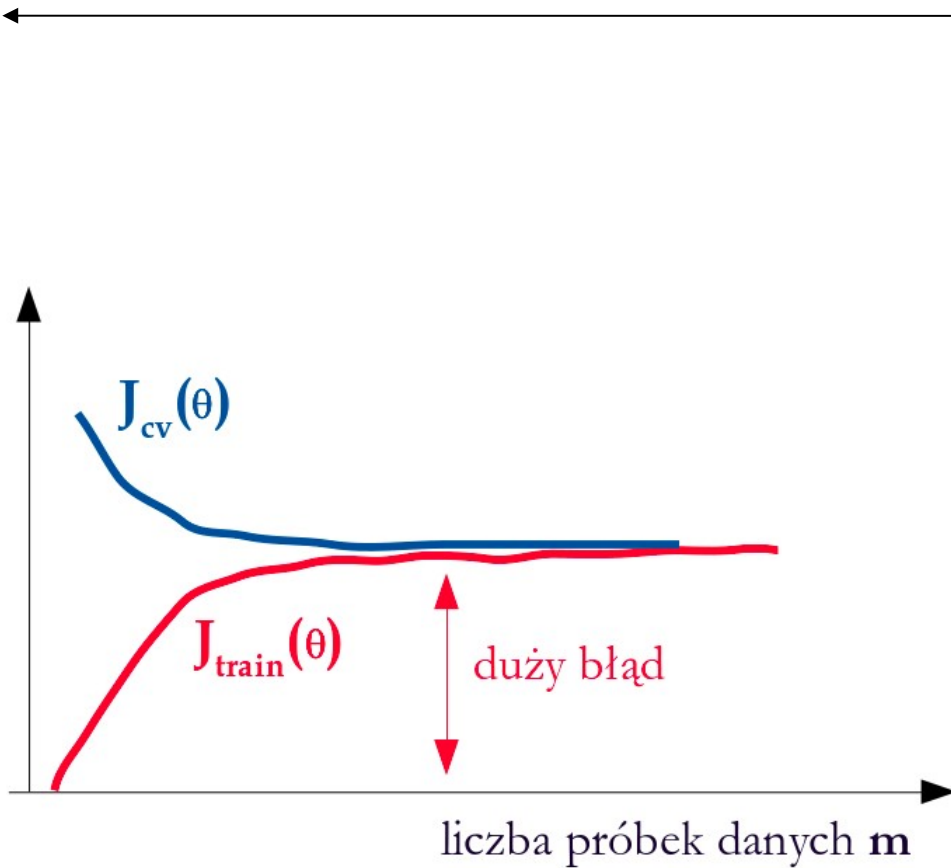
# Krzywe uczenia

*learning curves*

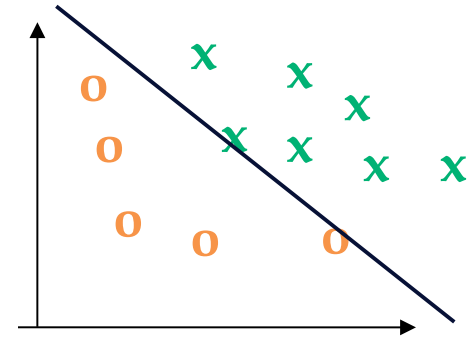
Wyznaczamy je dla konkretnego modelu  
(ustalony stopień wielomianu, ustalony parametr  $\lambda$ )  
sztucznie ograniczając próbkę danych:



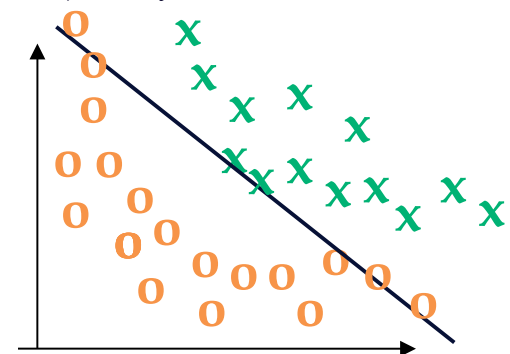
# Krzywe uczenia – przykład high bias



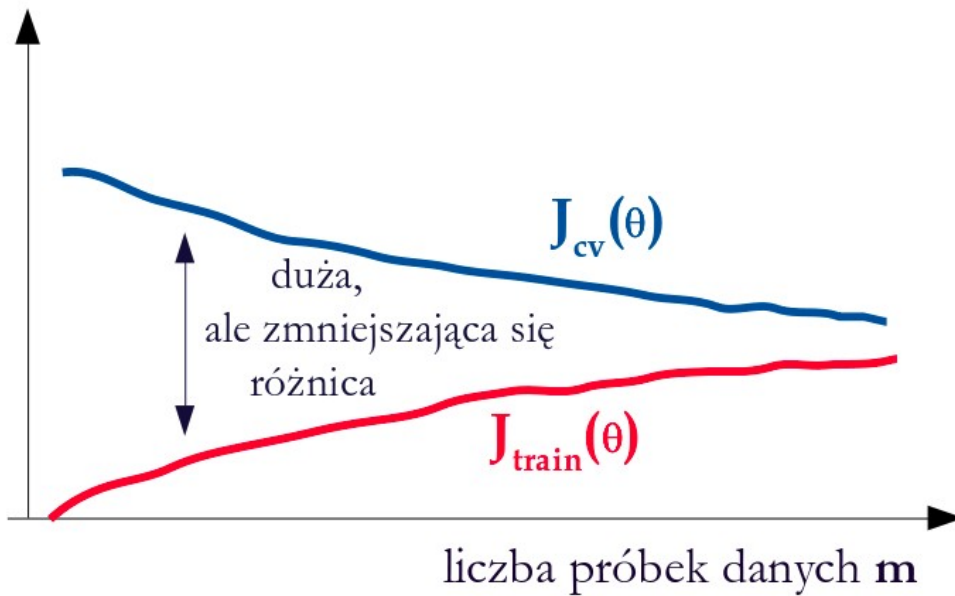
mało danych:



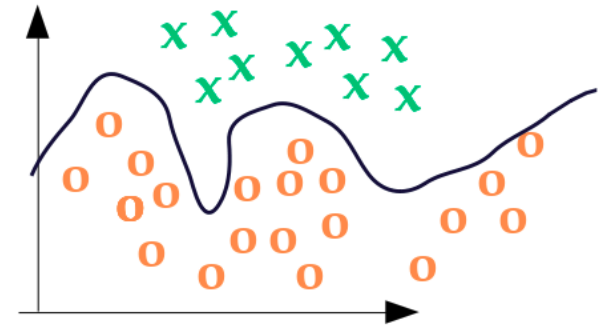
więcej danych:



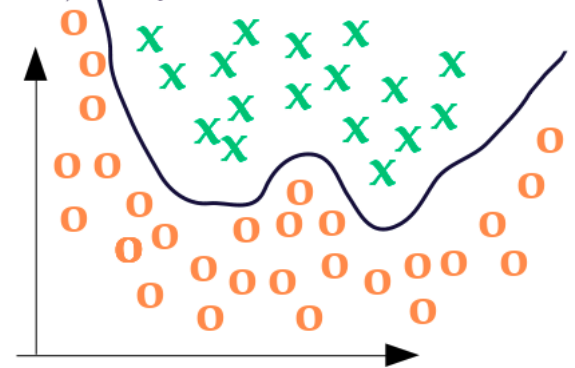
# Krzywe uczenia – przykład high variance



mało danych:

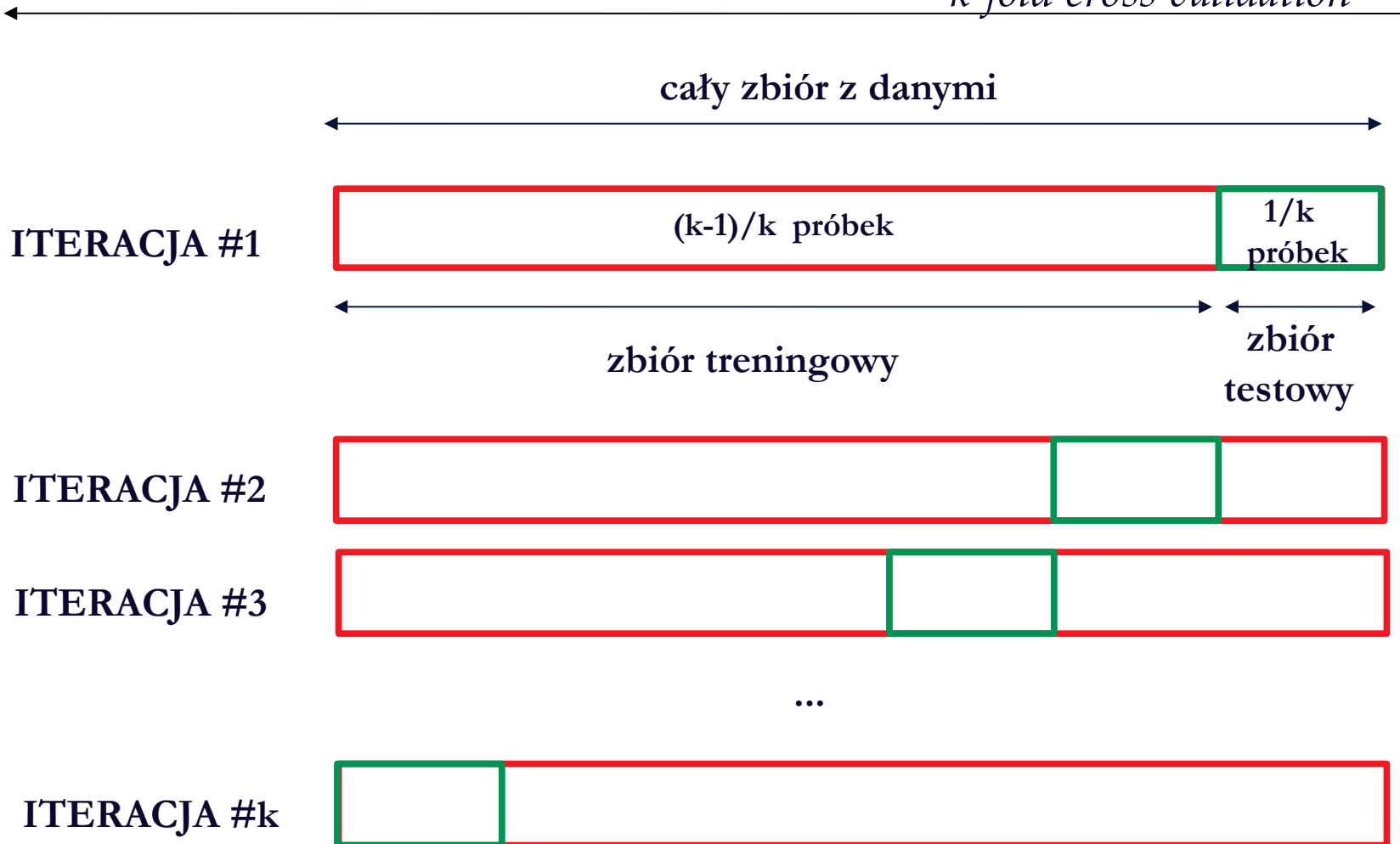


więcej danych:



# Walidacja krosowa

*k-fold cross validation*



⇒ końcowe parametry (np. *accuracy*) uśredniamy po  $k$  iteracjach

# Diagnostyka - podsumowanie

---

## Możliwe strategie:

- dodać więcej danych wejściowych (próbek)?  
→ pomoże w przypadku *high variance*
- zwiększyć wymiar danych (liczbę cech)?  
→ pomoże w przypadku *high bias*
- dodać nowe dane – wielomiany cech istniejących?  
→ pomoże w przypadku *high bias*
- zmniejszyć wymiar danych?  
→ pomoże w przypadku *high variance*
- regularyzacja – większa  $\lambda$ ?  
→ pomoże w przypadku *high variance*
- regularyzacja – mniejsza  $\lambda$ ?  
→ pomoże w przypadku *high bias*

# Prosty przykład kodowy \*

---

\* załączony plik jupyter-notebook:

<http://tele.agh.edu.pl/~kulakowski/ml/logistic-regression.ipynb>

**Dziękuję za spotkanie!**

**Pytania?**

**Komentarze?**