

МАШИННОЕ ОБУЧЕНИЕ И АНАЛИЗ ДАННЫХ

(Machine Learning and Data Mining)

Н. Ю. Золотых

<http://www.uic.unn.ru/~zny/ml>

Лекция 17

Дilemma «Смещение–разброс» и кривая обучения

Bias-variance trade-off и Learning curve

17.1. Дileмма «Смещение–разброс»

$$Y = f^*(x) + \varepsilon(x), \quad \mathbb{E} \varepsilon(x) = 0, \quad \text{Var } \varepsilon(x) = \sigma_\varepsilon^2(x)$$

$f(x, D)$ – решающее правило,

$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$ – обучающая выборка

$x, f^*(x)$ случайными не являются

Случайными являются: $\varepsilon(x)$ и D , причем считаем их независимыми

Тогда

$$\text{MSE}(x) = \mathbb{E}_{D, \varepsilon} (f(x, D) - Y)^2 = \sigma_\varepsilon^2(x) + \underbrace{\mathbb{E}_D^2 (f(x, D) - f^*(x))}_{\text{Bias}^2 x} + \underbrace{\mathbb{E}_D (f(x, D) - \mathbb{E}_D f(x, D))^2}_{\text{Variance } x}$$

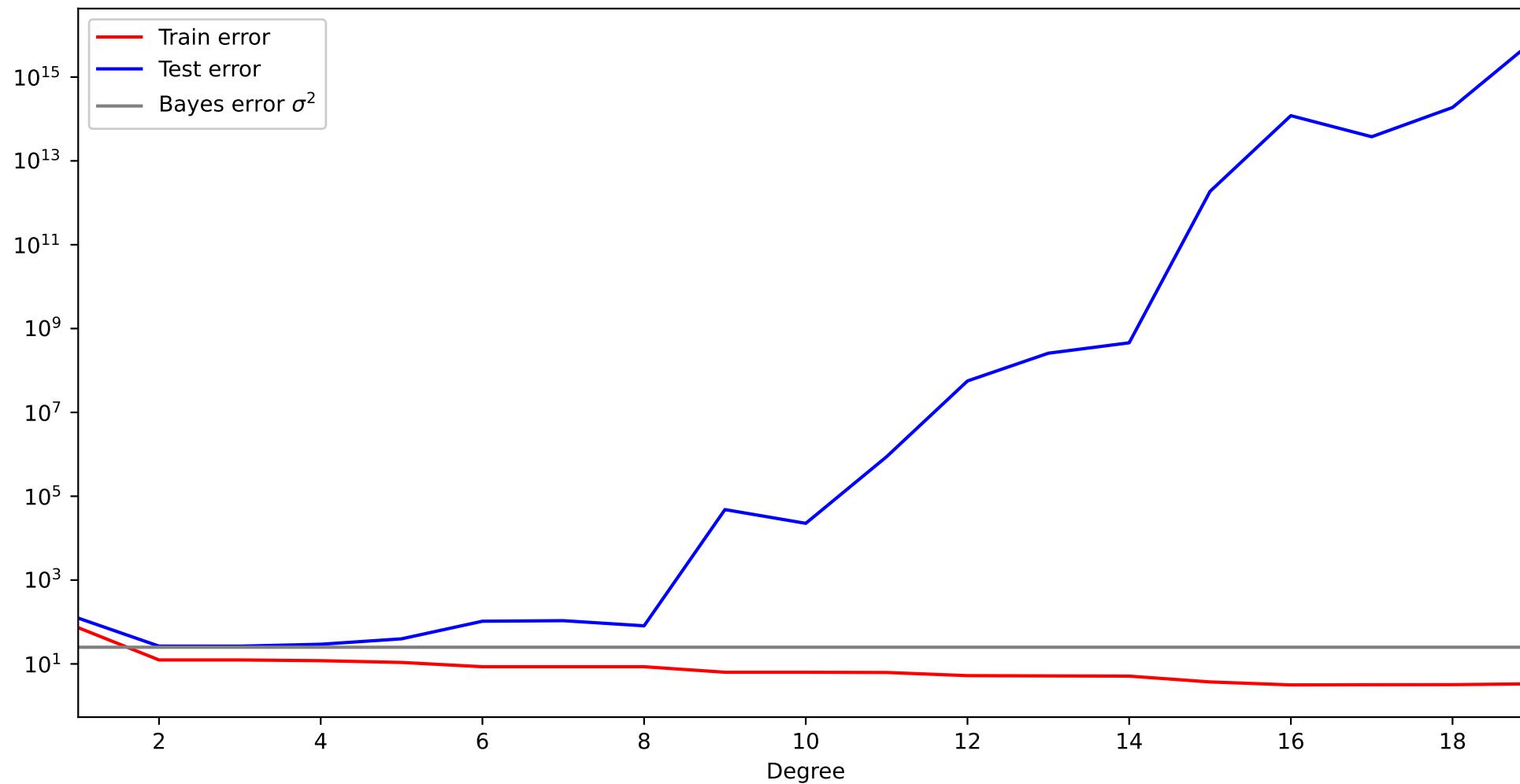
Ошибка = Неустранимая ошибка + Смещение² + Разброс

Доказательство:

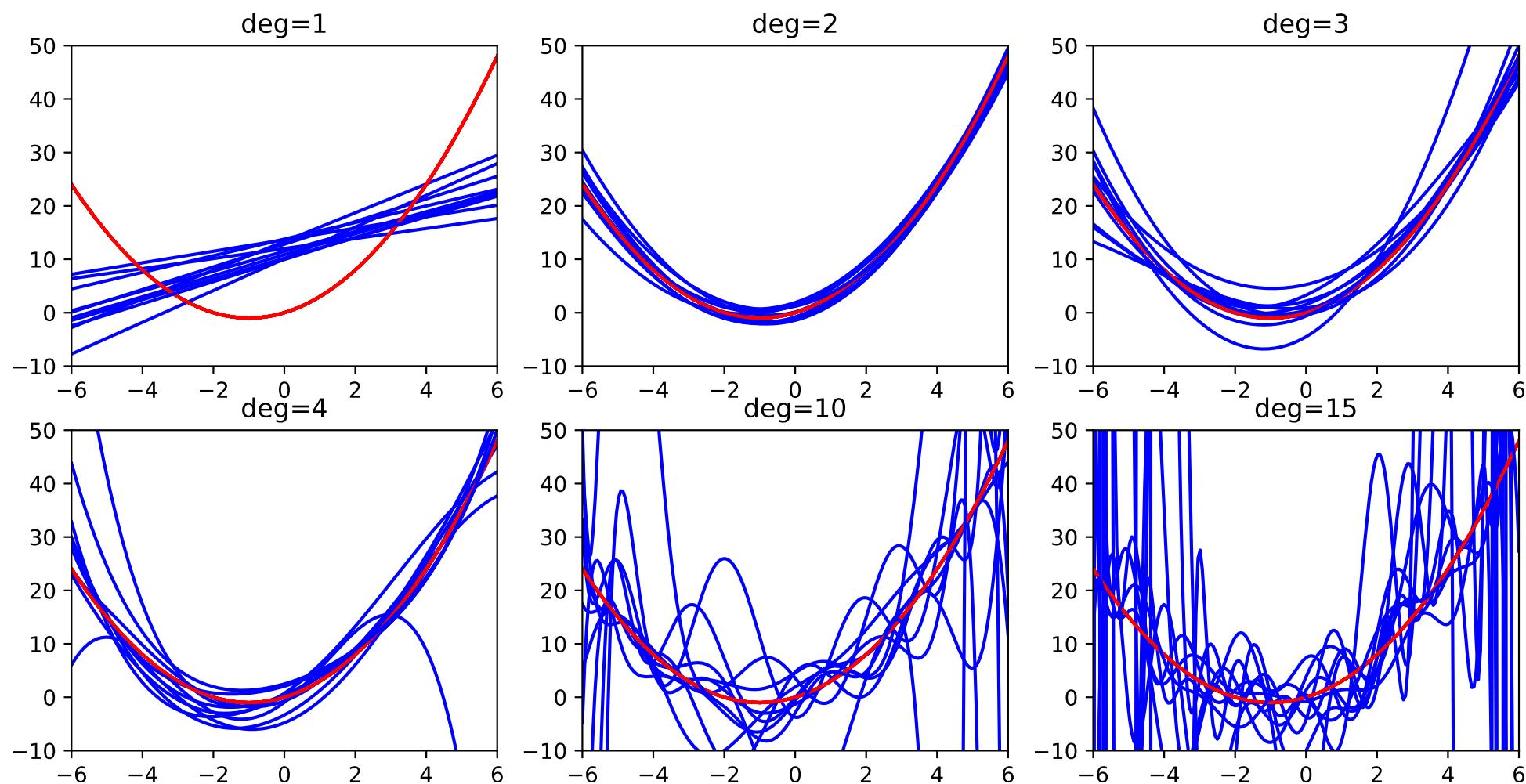
(можно раскрыть л.ч и п.ч. и увидеть, что они совпадают)

Большое смещение – недообучение. Большой разброс – переобучение.

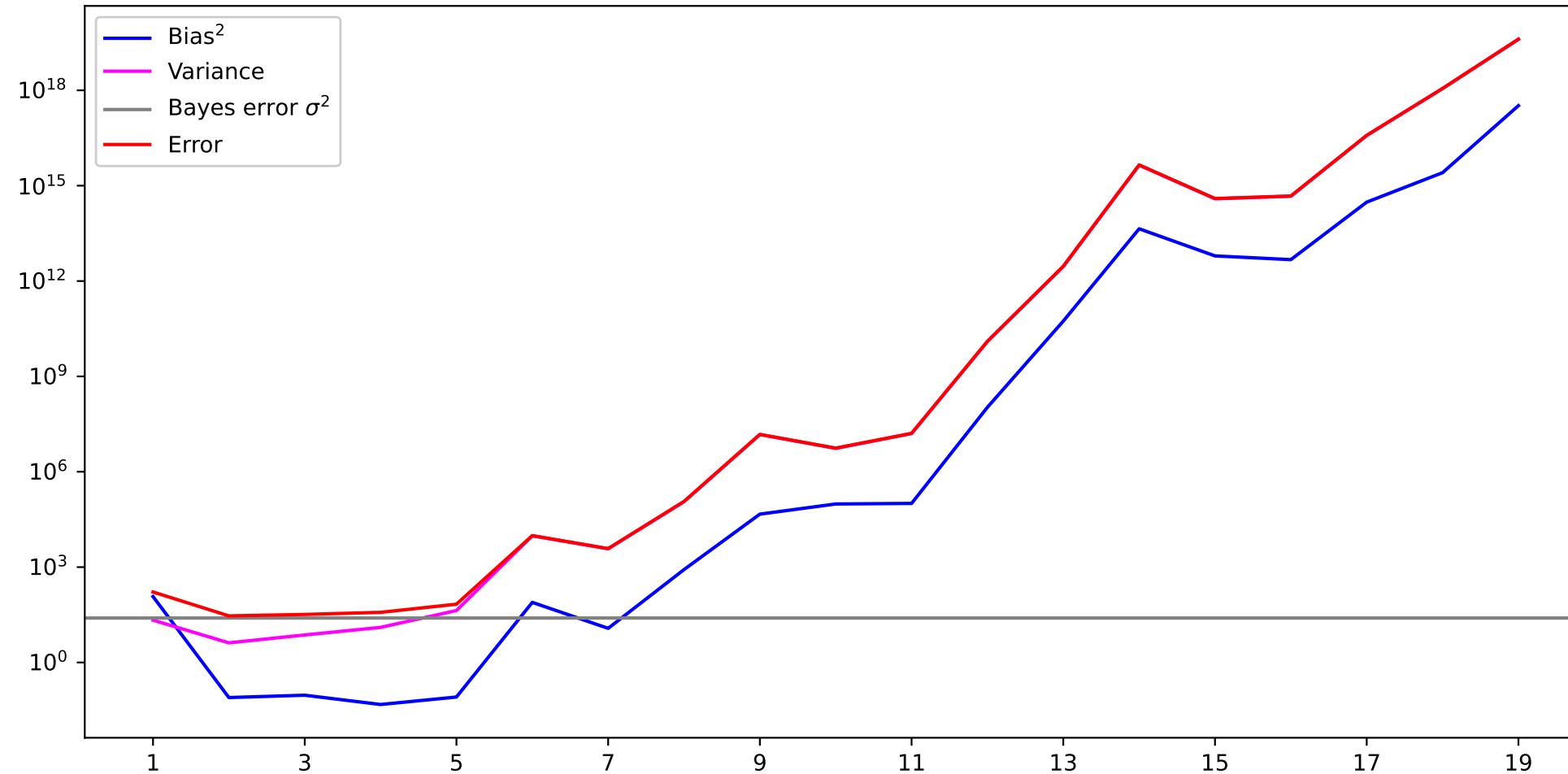
LSQ $y = x^2 + 2x + N(0, 5)$



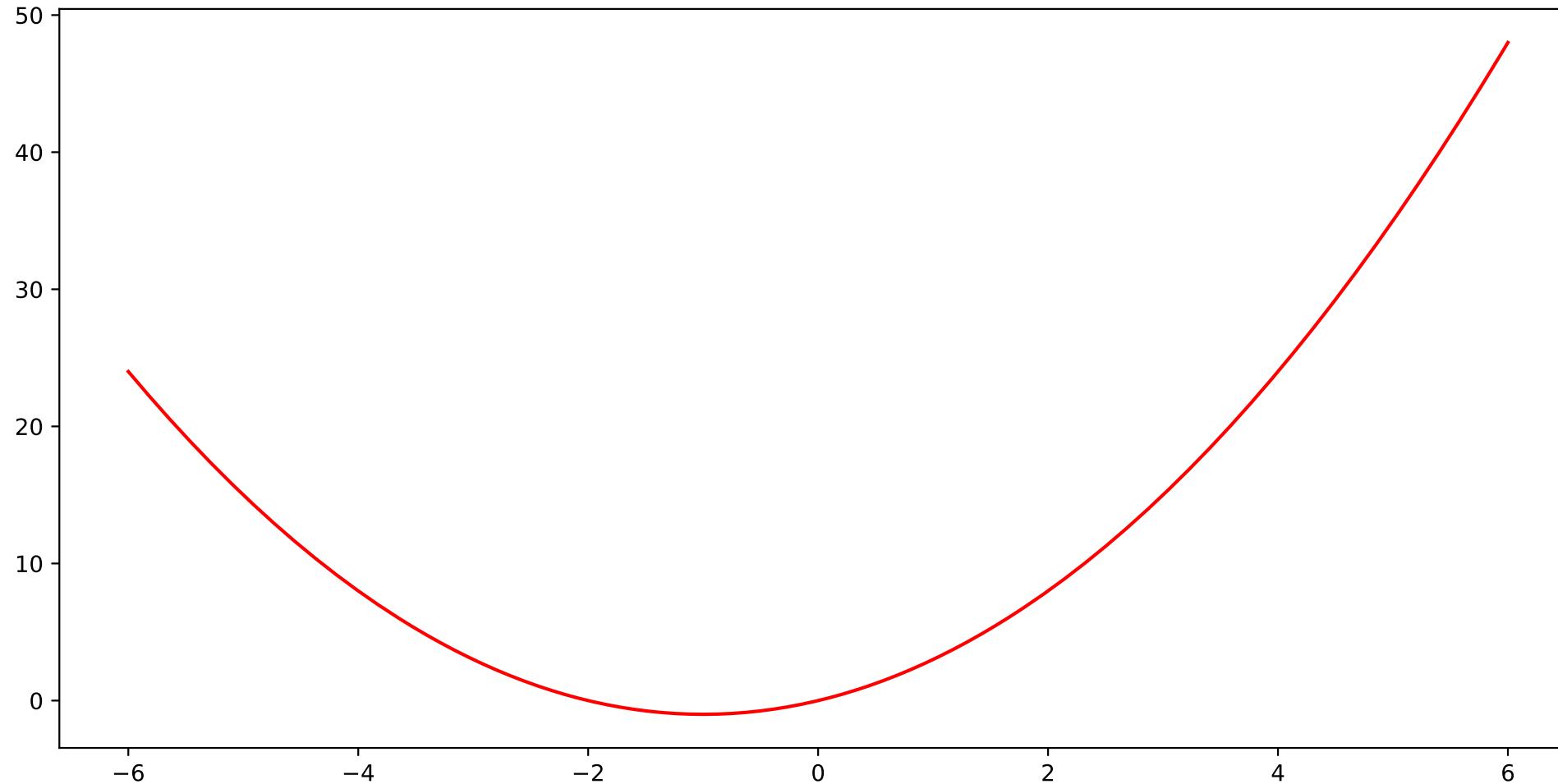
LSQ



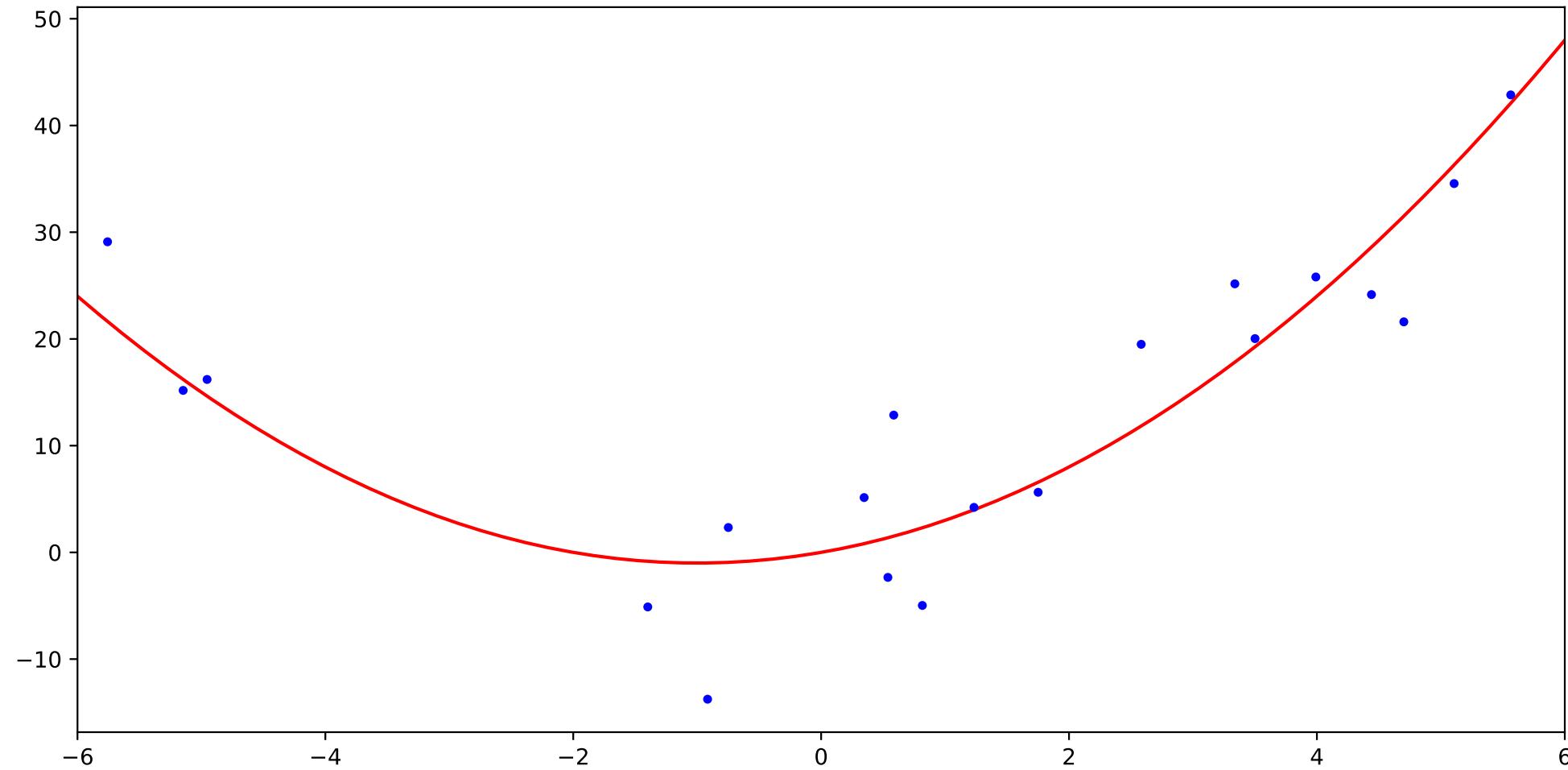
LSQ



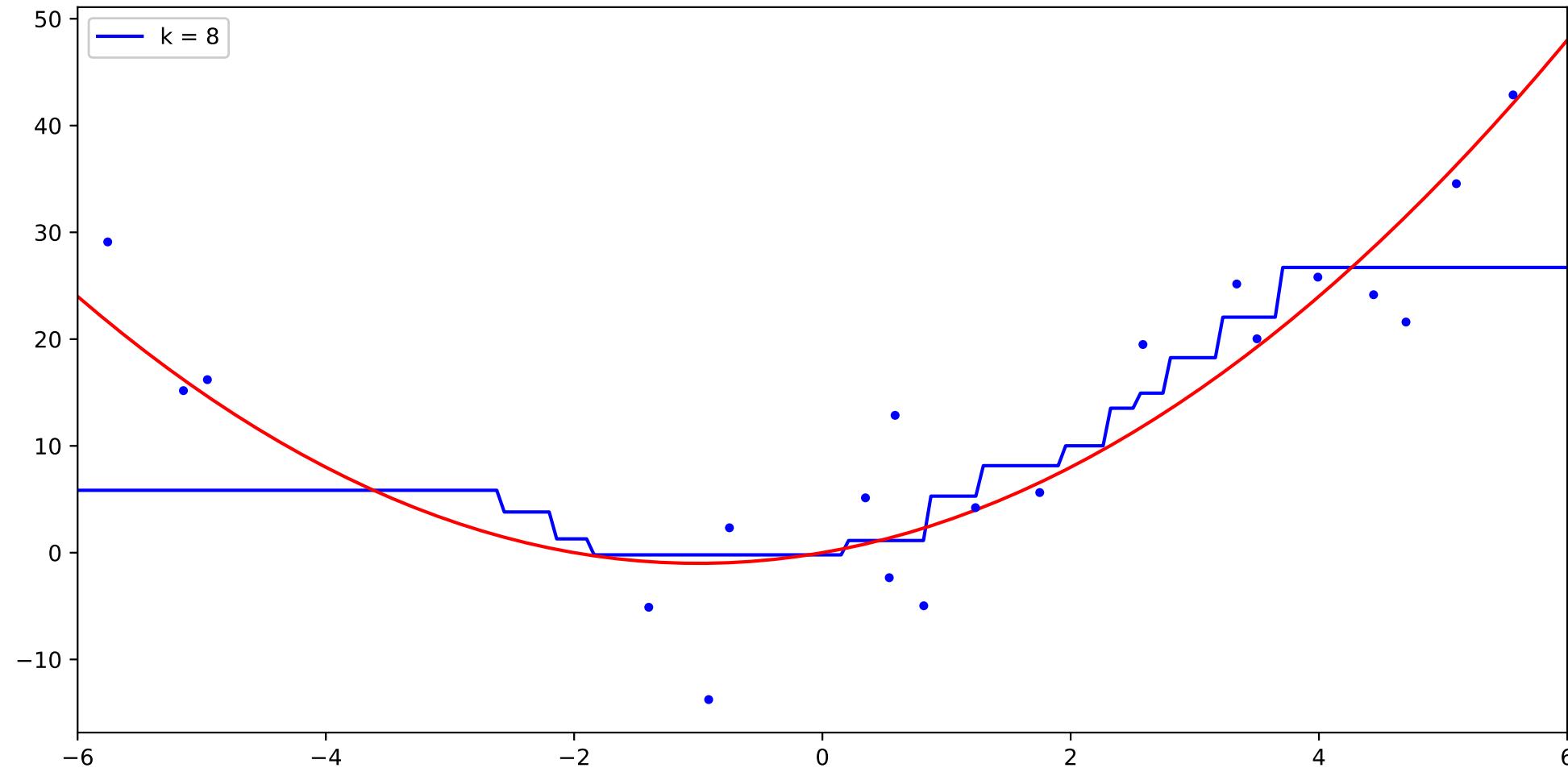
kNN $y = x^2 + 2x + N(0, 5)$



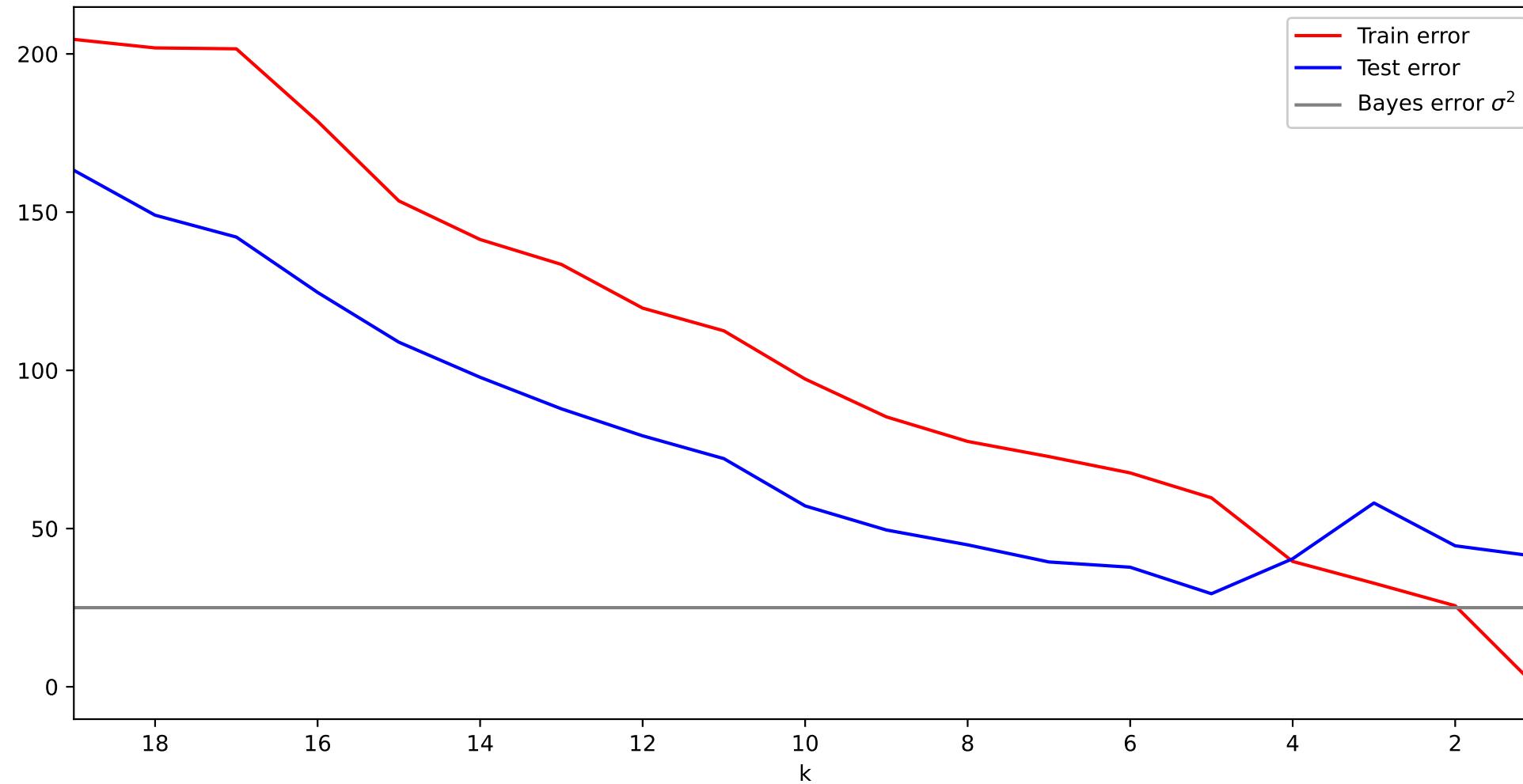
kNN $y = x^2 + 2x + N(0, 5)$

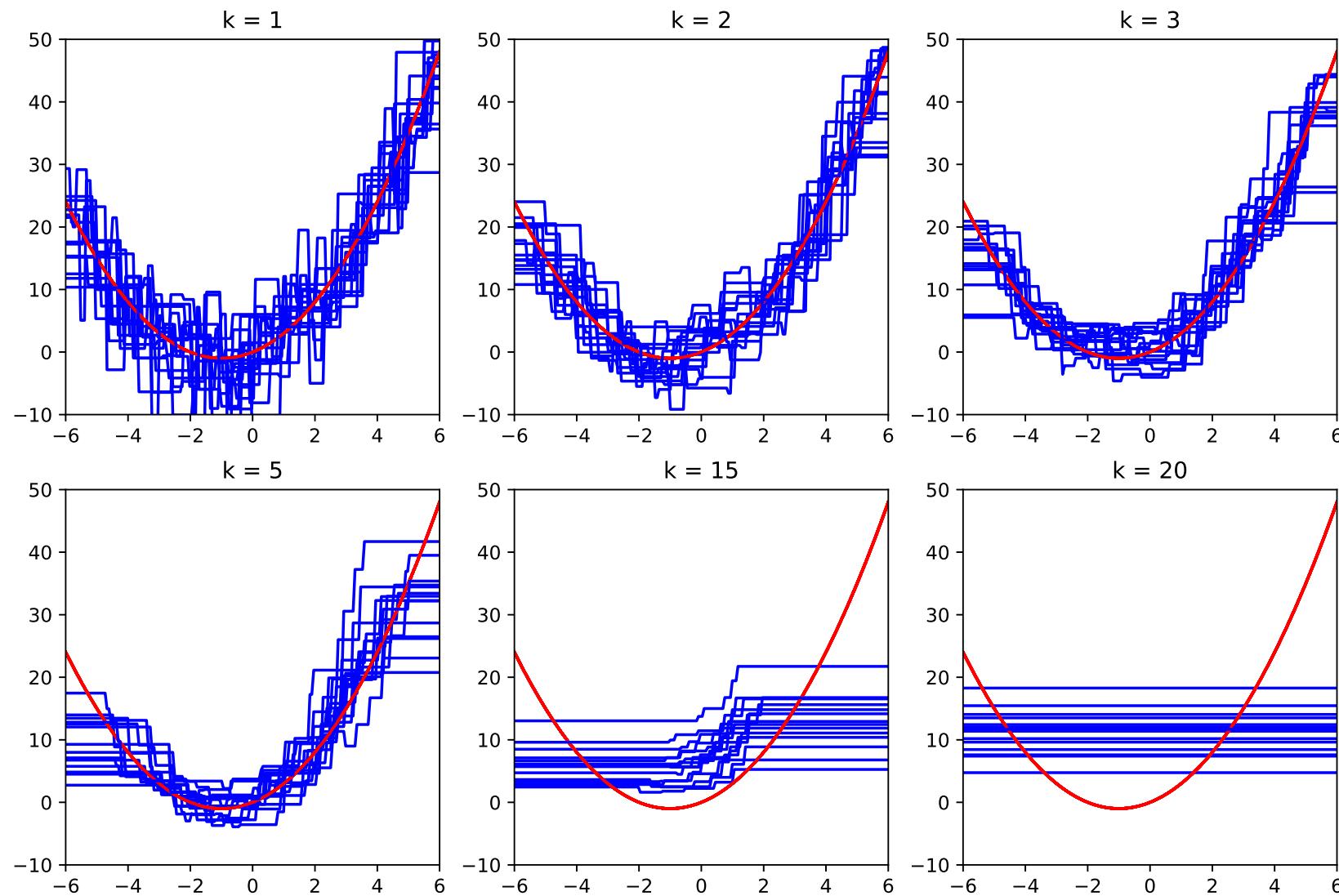


$$\text{kNN } y = x^2 + 2x + N(0, 5)$$

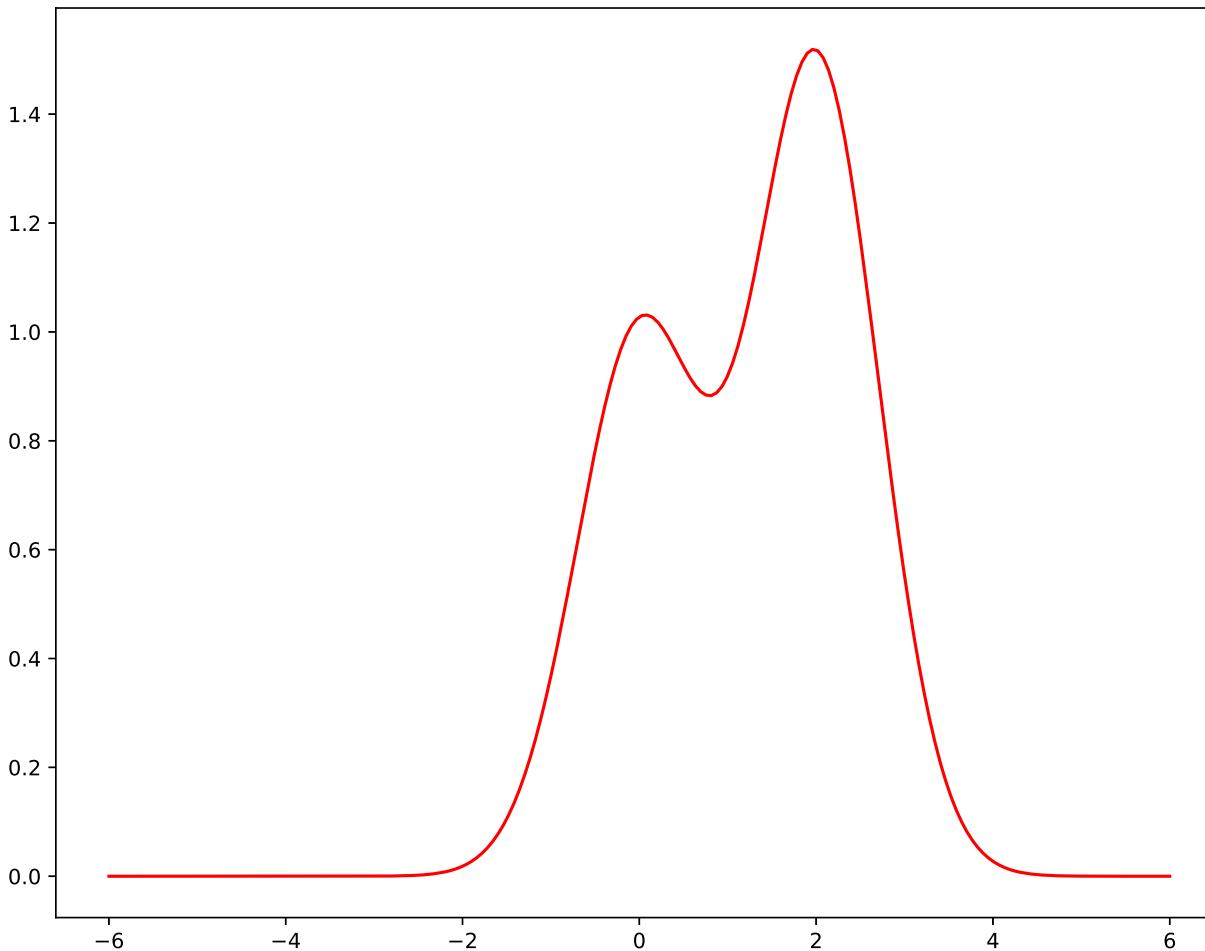


kNN $y = x^2 + 2x + N(0, 5)$

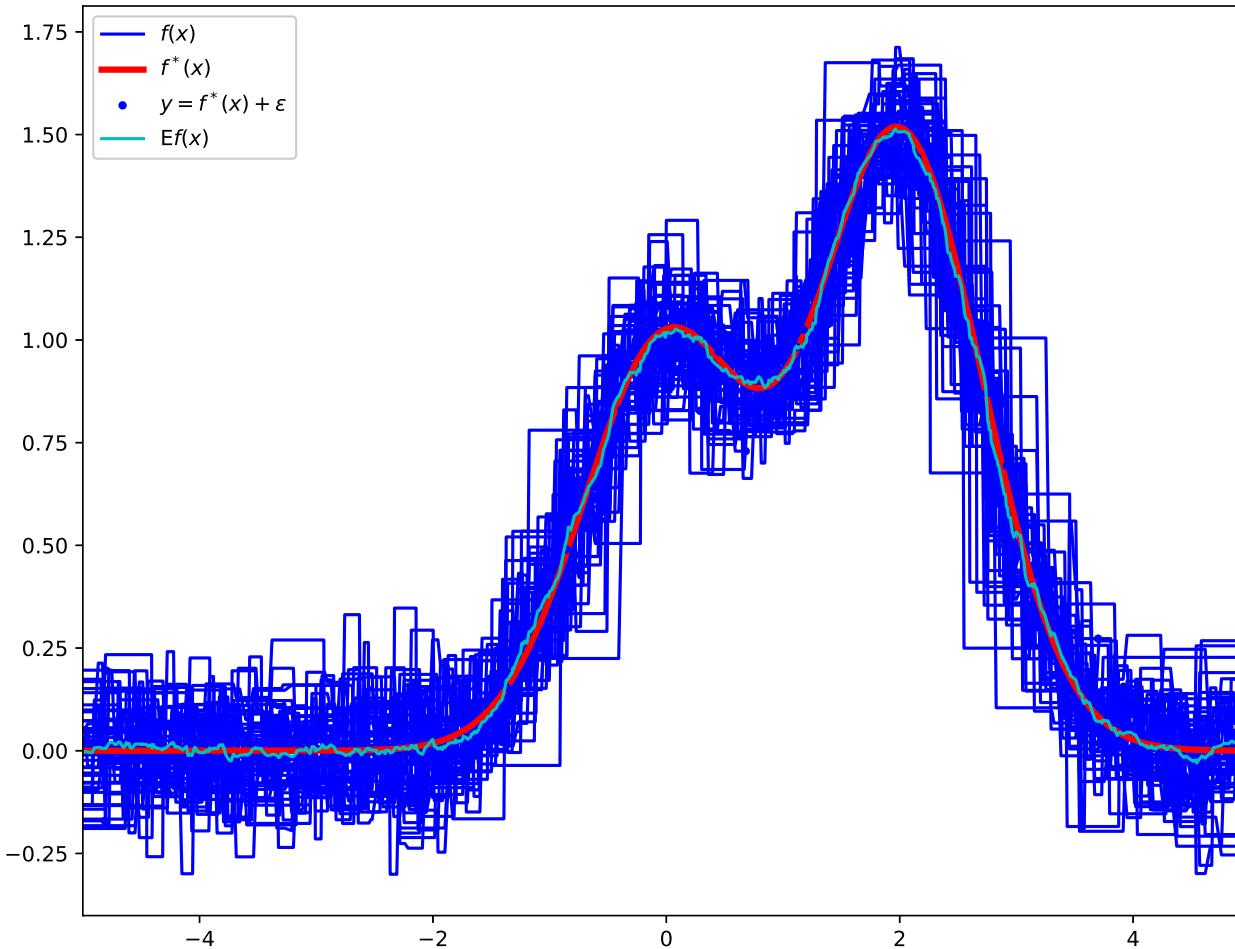




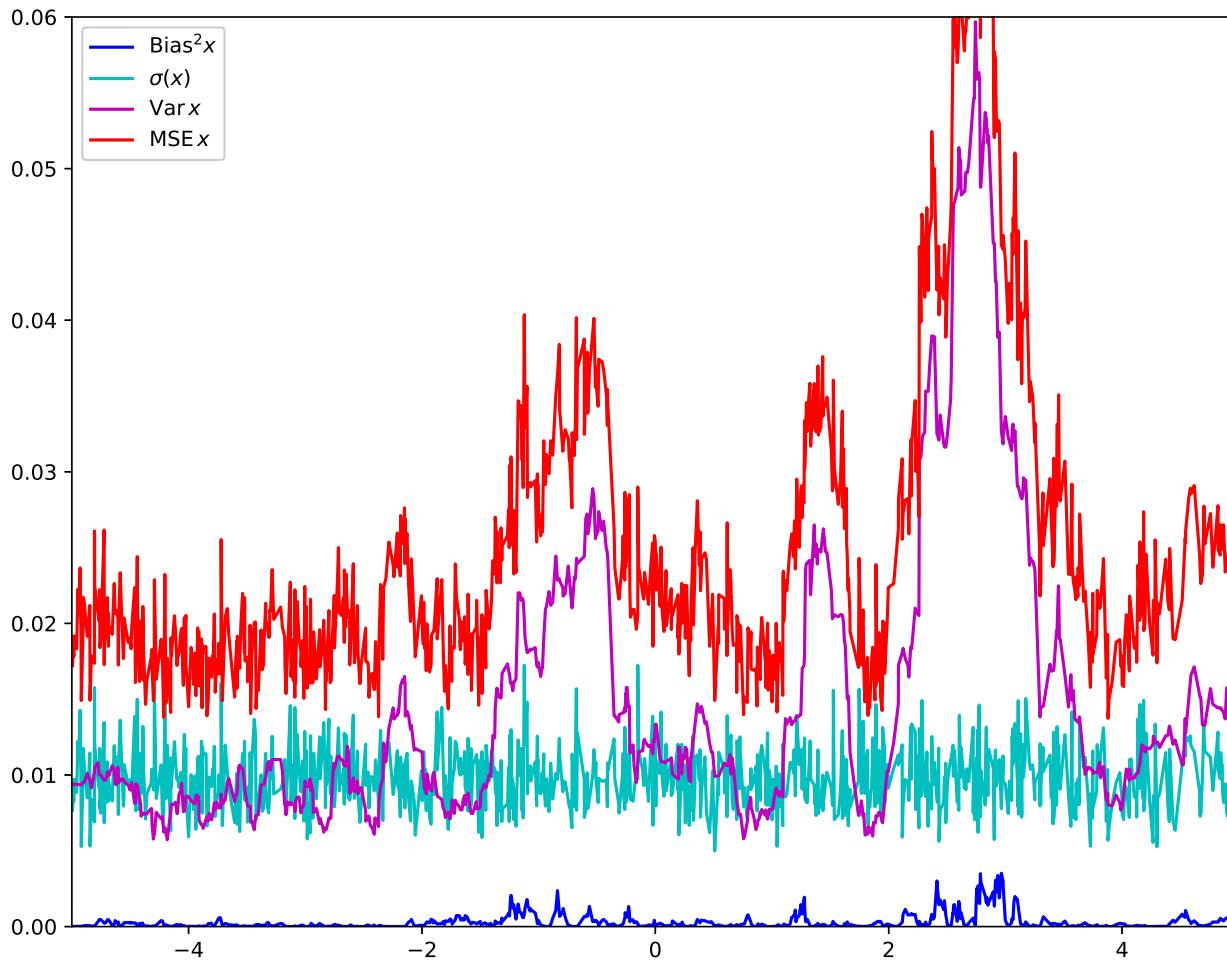
$$y = e^{-x^2} + 1.5e^{-(x-2)^2}$$



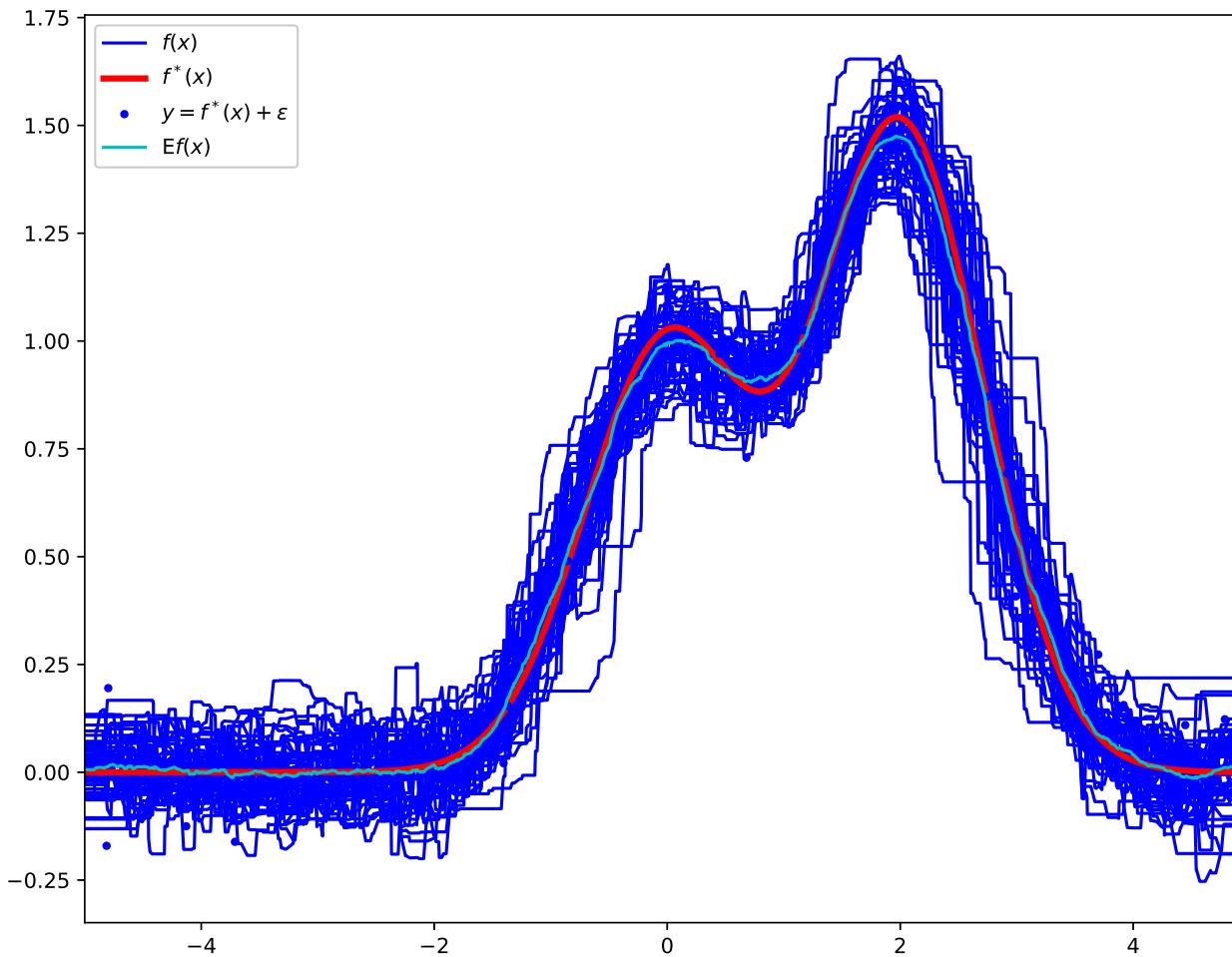
$$\text{DT } y = e^{-x^2} + 1.5e^{-(x-2)^2}$$



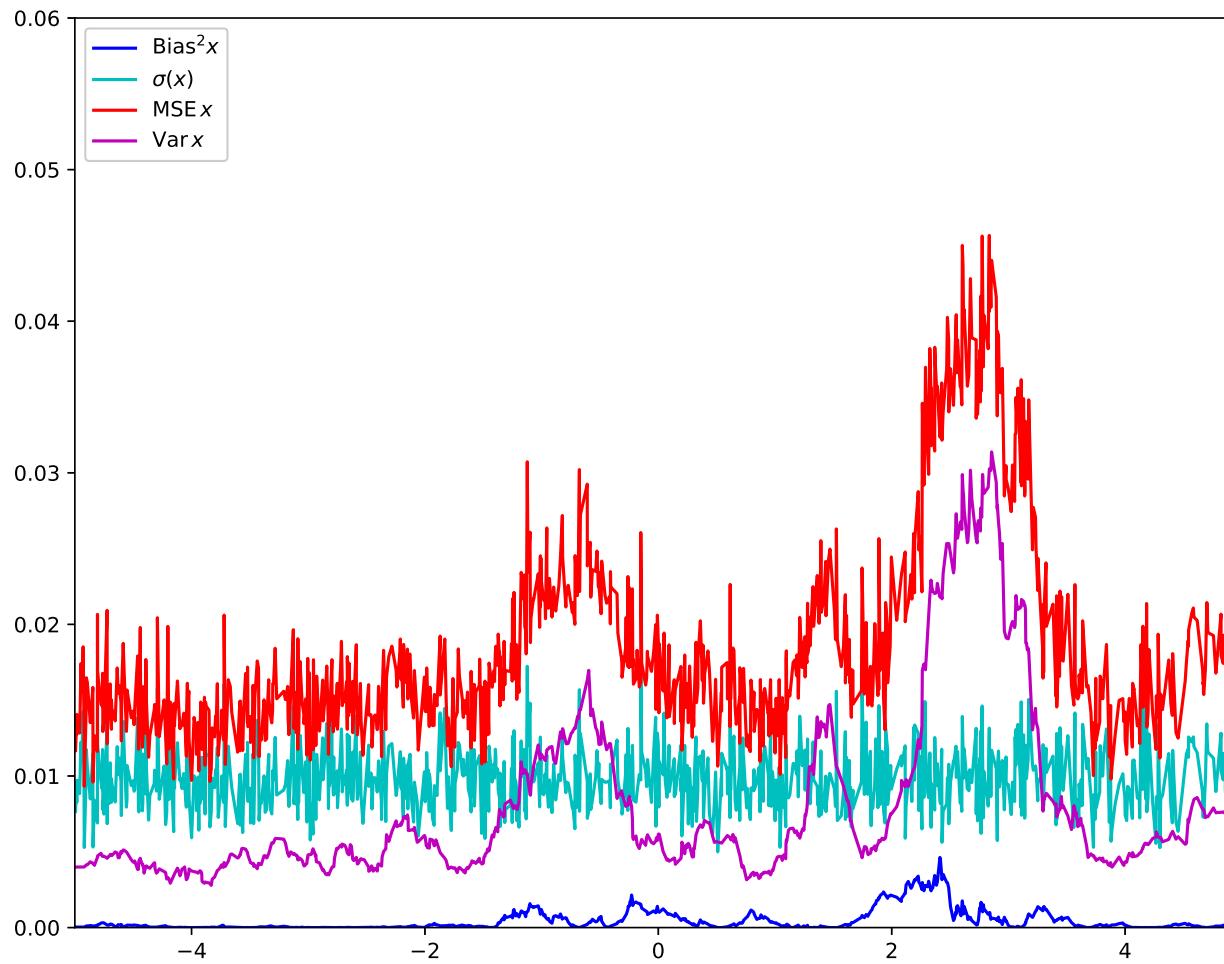
$$\text{DT } y = e^{-x^2} + 1.5e^{-(x-2)^2}$$



$$\text{RF } y = e^{-x^2} + 1.5e^{-(x-2)^2}$$



$$\text{RF } y = e^{-x^2} + 1.5e^{-(x-2)^2}$$



17.2. Кривая обучения (learning curve)