# Linear Regression

Linear regression . fit

$$\hat{\beta} = \underbrace{(X^T X)^{-1}}_{\text{1. matrix multiply}} X^T y \Big\} \text{ matrix multiply}$$

2. matrix solve

1. Matrix multiply ($X$ is $n \times p$)

$$(X^T X)_{jk} = \sum_i X_{ij} X_{ik} \quad - \quad O(n)$$

Computational complexity:

$$O(p^2 n)$$

2. Matrix solve — Cholesky decomp.

Computational complexity:

$$O(p^3)$$

Total Regression Fit Complexity:

$$O(p^3 + p^2 n)$$

Linear Regression, predict

$$\hat{f}(x^*) = x^{*T}\hat{\beta}$$

$O(p)$ time

▷ predict is fast $O(p)$
    fit is slow $O(p^3 + p^2 n)$

▷ $X^T X$ needs to be invertible
    (impossible if $p > n$!)

▷ $\hat{f}$ is linear

▷ Gauss - Markov theorem

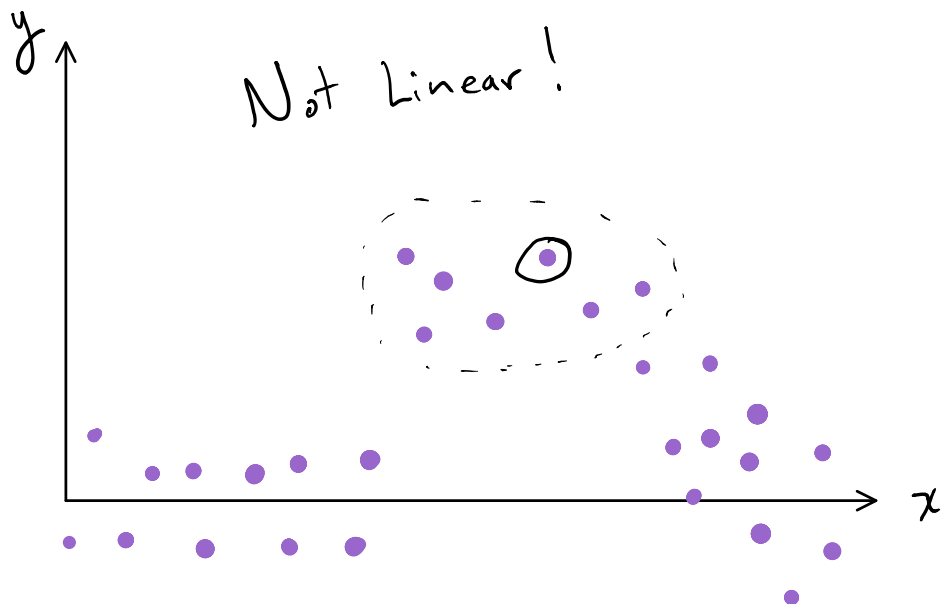if $y_i = x_i^T \beta + \varepsilon_i$ with $\mathbb{E}\varepsilon_i = 0$

$\mathbb{V}\varepsilon_i = \sigma^2$ then OLS is

unbiased: $\mathbb{E}\hat{\beta} = \beta$

and is the unbiased estimator

with minimum variance.

Not Linear!

Given a metric $d(x, x')$ then
the $k$-nearest neighbors of $x$ in
$\{x_i\}_{i=1}^n$ is $x_{m_1}, \cdots, x_{m_k}$ s.t.

$$d(x_{m_1}, x) \leq d(x_{m_2}, x) \leq \cdots \leq d(x_{m_n}, x)$$

<u>kNN methods</u> fit a model for
predicting $x$ using the kNN
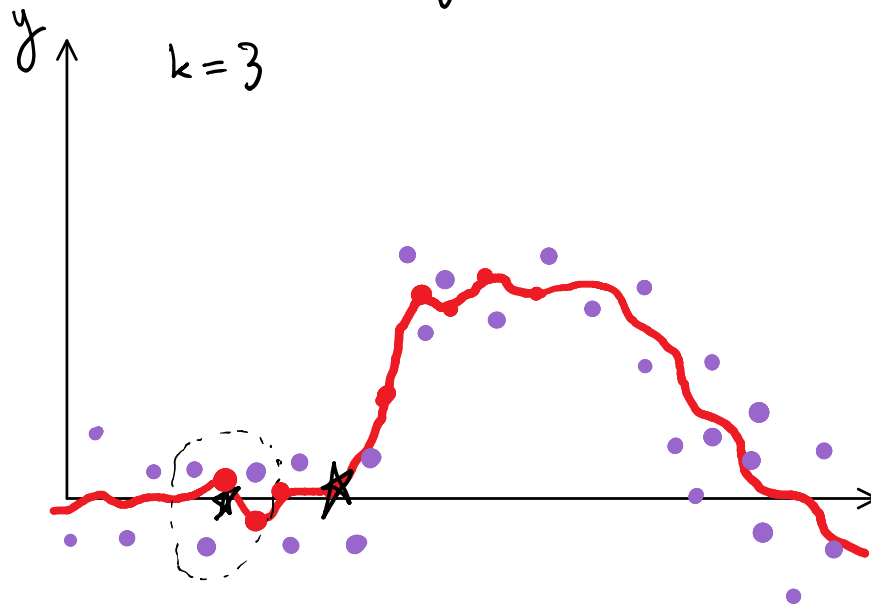dataset $\{x_{m_j}, y_{m_j}\}_{j=1}^k$.

Regression: $\hat{f}(x) = \frac{1}{k} \sum_{j=1}^k y_{m_j}$

Classification: $\hat{f}(x) = \mathbb{1}\left\{ \frac{1}{k} \sum_{j=1}^k y_{m_j} > \frac{1}{2} \right\}$
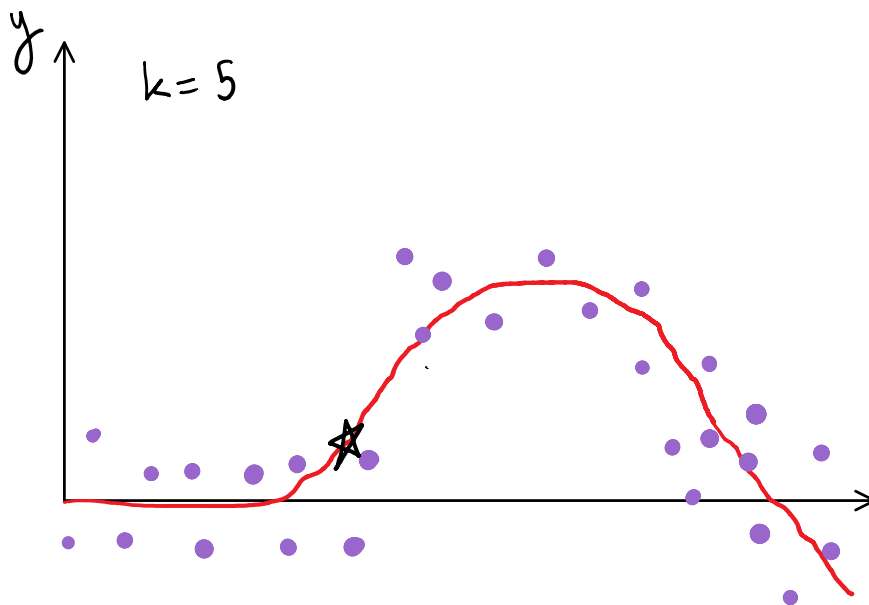
▷ randomized rounding if $= \frac{1}{2}$

$y$

Variance tracking $k = 2$

$y$

$k = 3$

low bias

high variance

$y$

$k = 5$

high bias

low variance

$$\{y_i, x_i\} \ iid$$

$$\mathbb{E}\left[\frac{1}{k}\sum_{j=1}^{k} y_{m_j} \Big| x\right] = \frac{1}{k}\sum_{j=1}^{k} \mathbb{E}[y|x_{m_j}]$$

$$\mathbb{V}\left[\frac{1}{k}\sum_{j=1}^{k} y_{m_j} \Big| x\right] = \frac{1}{k^2}\sum_{j=1}^{k} \mathbb{V}[y|x_{m_j}] = \frac{1}{k}\sigma^2$$

Bias: $\mathbb{E}[\hat{f}(x)] - \mathbb{E}[y|x] = \frac{1}{k}\sum_{j=1}^{k} \mathbb{E}[y|x_{m_j}] - \mathbb{E}[y|x]$

as $\quad d(x_{m_j}, x) \nearrow \quad |\mathbb{E}[y|x_{m_j}] - \mathbb{E}[y|x]| \nearrow$

So, typically as $k \nearrow$ bias $\nearrow$

and

as $k \nearrow$ Variance $\searrow$

this is the bias-variance tradeoff
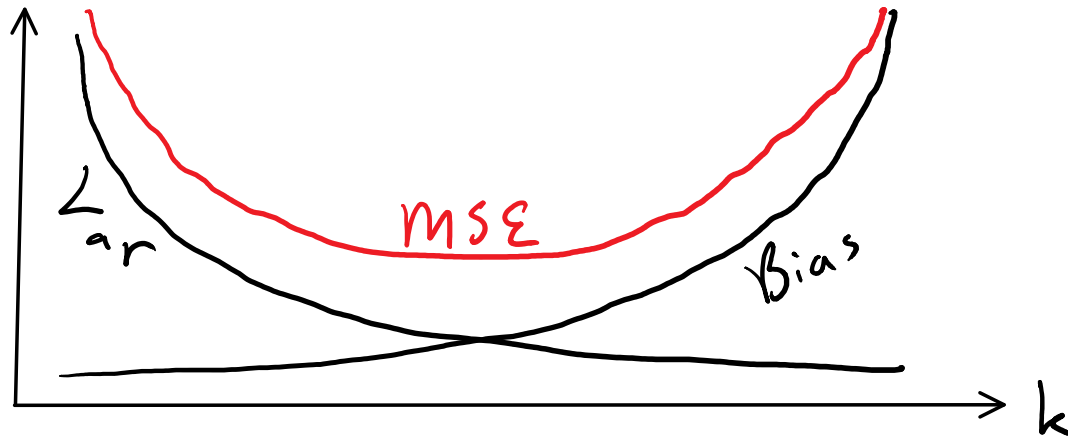
$$MSE = \mathbb{E}(\hat{f}(x) - \mathbb{E}[y|x])^2 \quad (\text{True Risk})$$

$$= \mathbb{E}(\hat{f}(x) - \mathbb{E}\hat{f}(x) + \mathbb{E}\hat{f}(x) - \mathbb{E}[y|x])^2$$

$$= \mathbb{E}(\hat{f}(x) - \mathbb{E}\hat{f}(x))^2 + (\mathbb{E}\hat{f}(x) - \mathbb{E}[y|x])^2$$
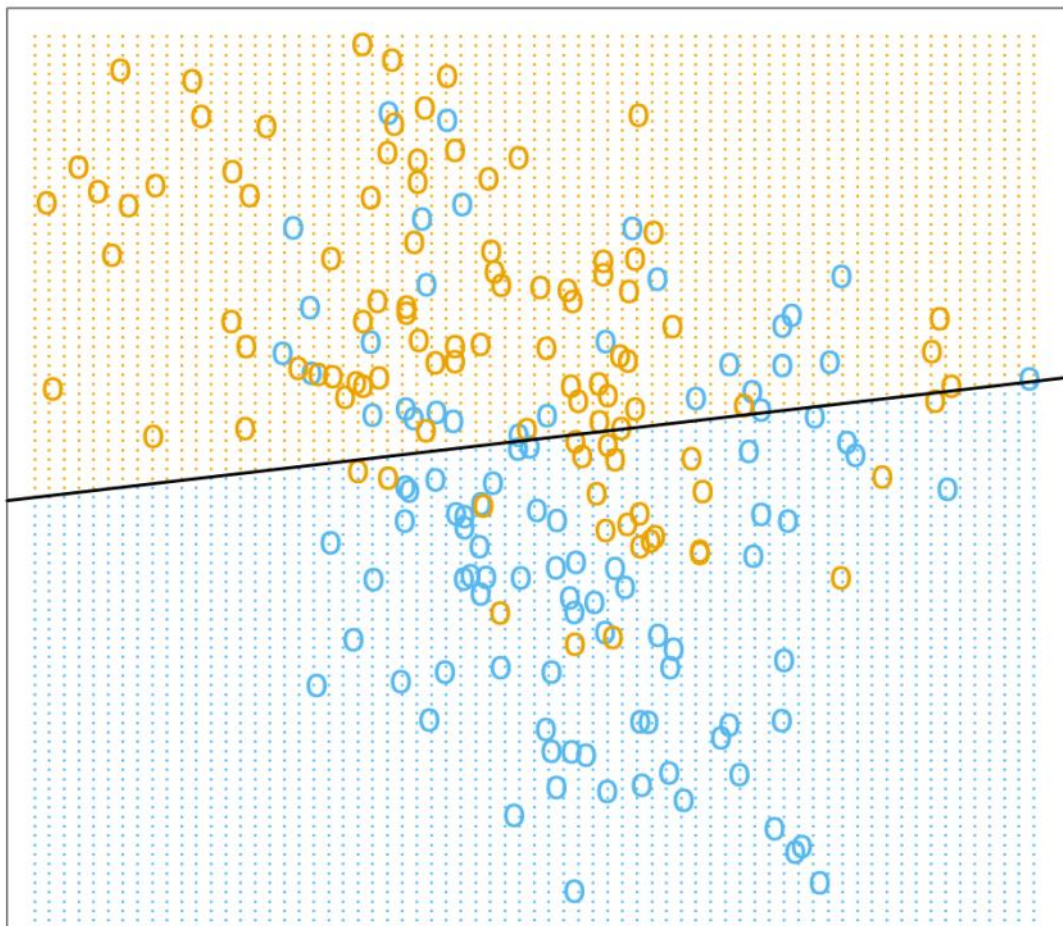
$$+ 2\,\mathbb{E}\underbrace{\left(\hat{f}(x) - \mathbb{E}\,\hat{f}(x)\right)}_{O} \cdot \left(\mathbb{E}\hat{f}(x) - \mathbb{E}[y|x]\right)$$

$$= bias^2 + \mathbb{V}\hat{f}(x)$$



ESL Ch. 2

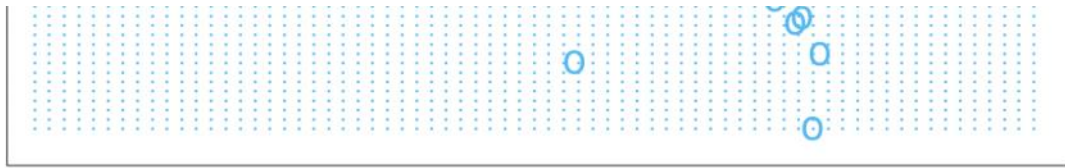Linear Regression of 0/1 Response

**FIGURE 2.1.** *A classification example in two dimensions. The classes are coded as a binary variable (*BLUE *= 0,* ORANGE *= 1), and then fit by linear regression. The line is the decision boundary defined by $x^T \hat{\beta} = 0.5$. The orange shaded region denotes that part of input space classified as* ORANGE, *while the blue region is classified as* BLUE.
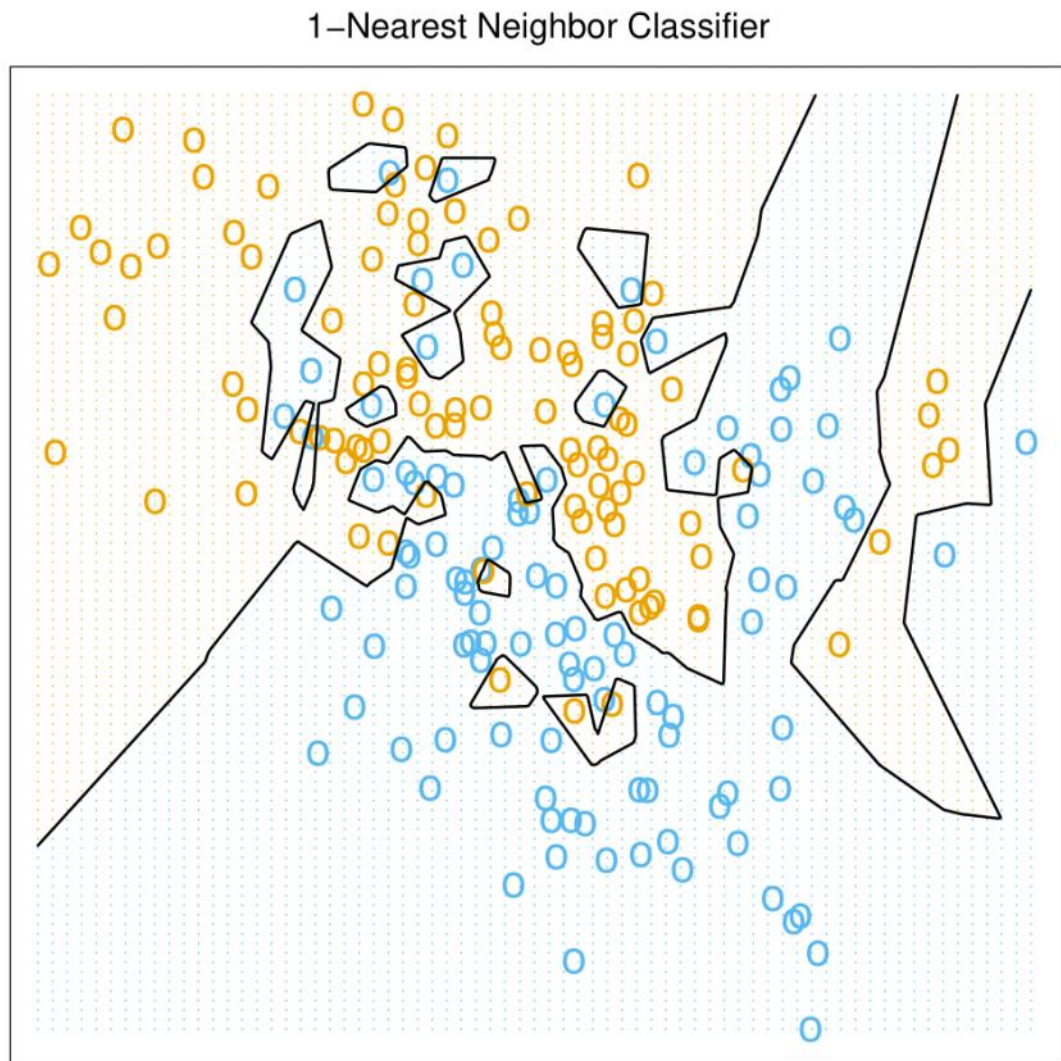
## 1–Nearest Neighbor Classifier



**FIGURE 2.3.** *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (*BLUE *= 0,* ORANGE *= 1), and then predicted by 1-nearest-neighbor classification.*
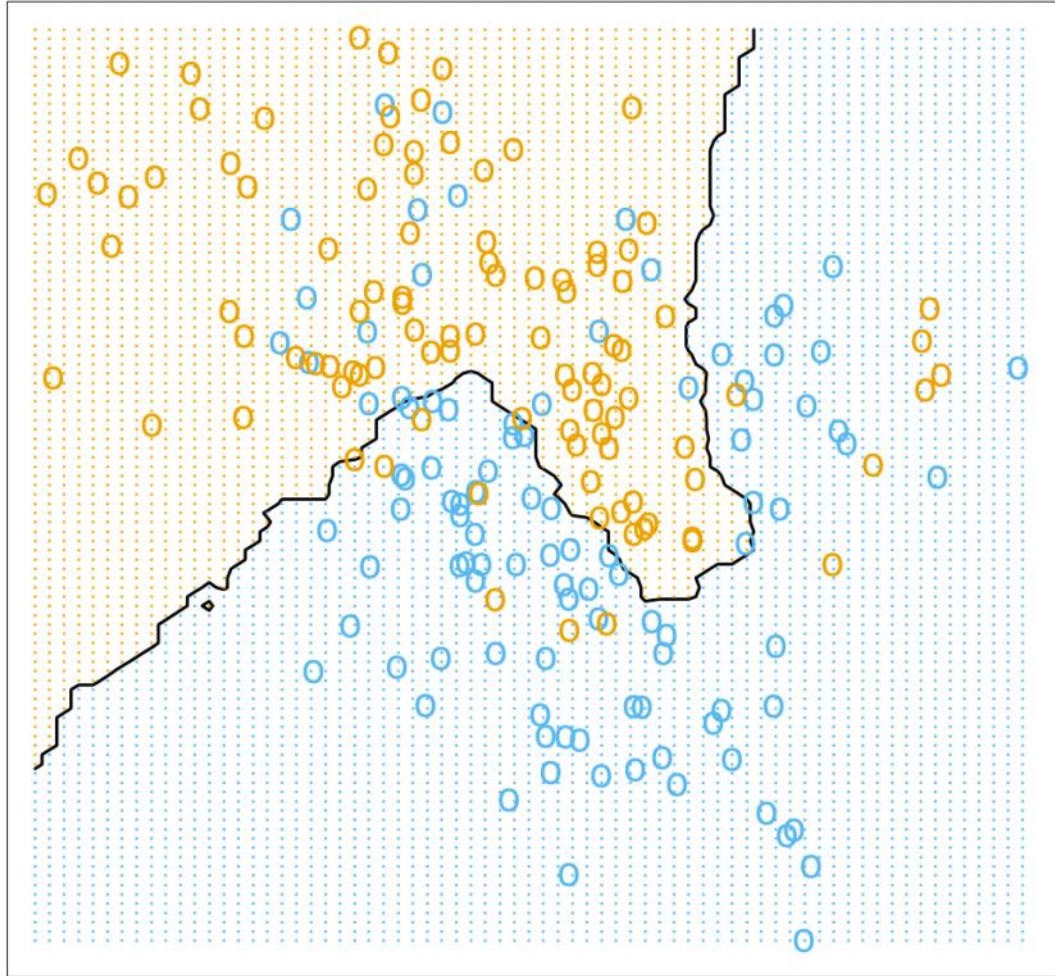
**FIGURE 2.2.** *The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable* (BLUE = 0, ORANGE = 1) *and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.*