# Machine Learning

FEM31002
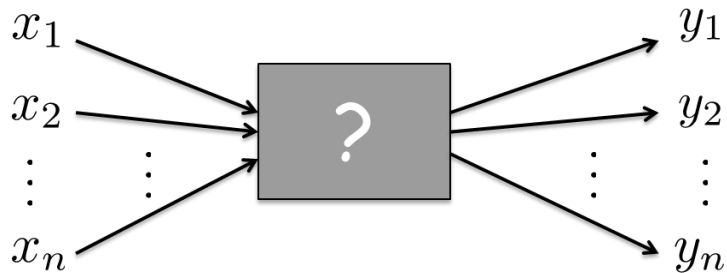
## Introduction

Part 2

**Ilker Birbil**

birbil@ese.eur.nl

$$Y = f(X) + \epsilon \quad \xrightarrow{\text{approximation?}} \quad \hat{Y} = \hat{f}(X)$$

$$x_1 \quad x_2 \quad \vdots \quad x_n \rightarrow \boxed{?} \rightarrow y_1 \quad y_2 \quad \vdots \quad y_n$$
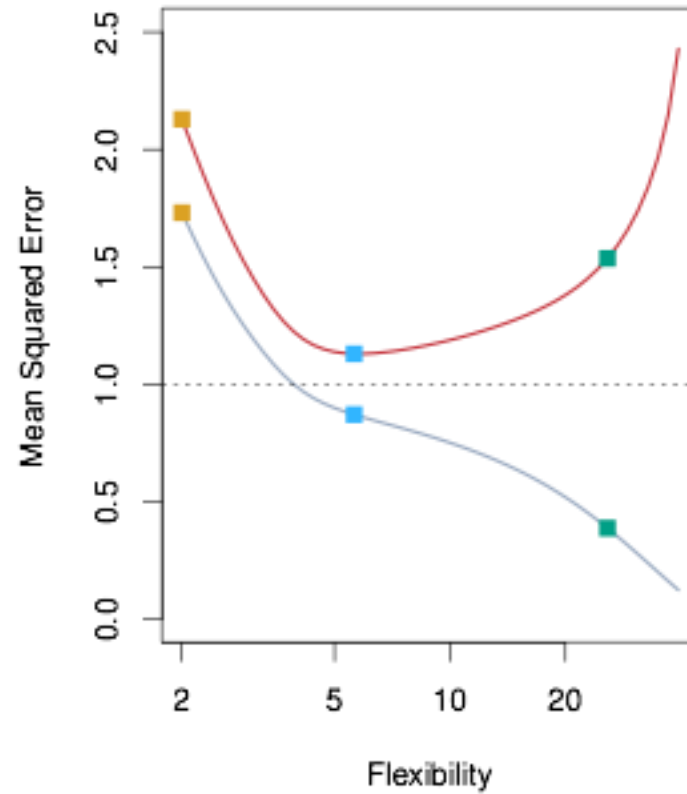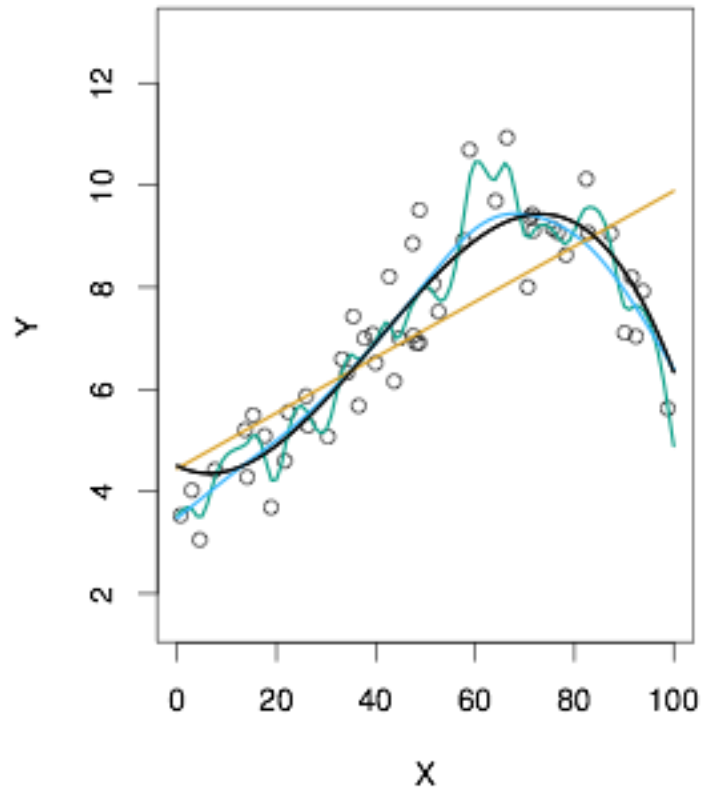
training data

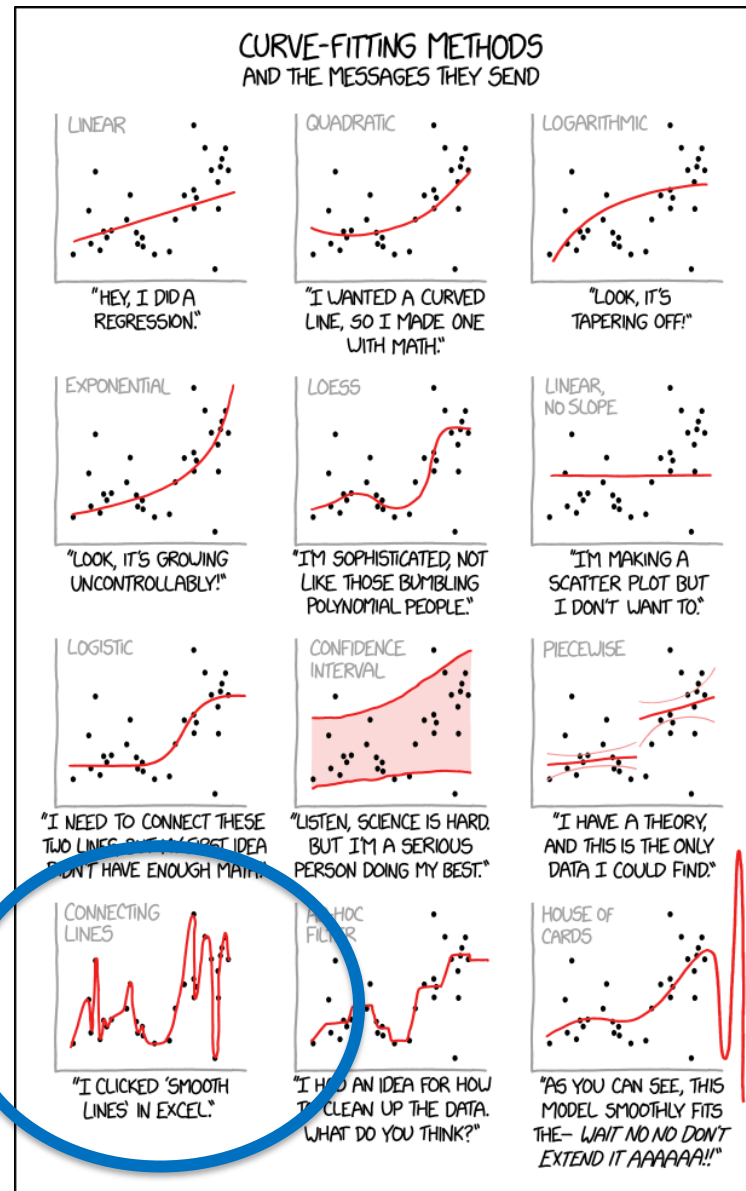$$\{(x_i, y_i) : 1, \ldots, n\}$$

**Training** Mean Square Error (MSE)

$$\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{f}(x_i))^2$$

$(x_0, y_0)$ : test data  ???

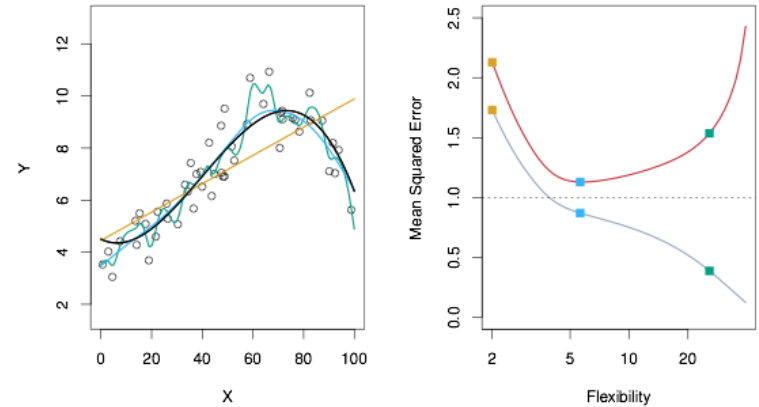$$\mathbb{E}(Y - \hat{Y})^2 = \left(f(X) - \hat{f}(X)\right)^2 + \mathrm{Var}(\epsilon)$$



overfitting

# Expected **test** MSE:

$$
\begin{aligned}
\mathbb{E}(y_0 - \hat{f}(x_0))^2 &= \mathbb{E}\left(\left(f(x_0) + \epsilon - \hat{f}(x_0)\right)^2\right) \\
&= \mathbb{E}\left((f(x_0))^2 + 2\epsilon f(x_0) + \epsilon^2 - 2(f(x_0) + \epsilon)\hat{f}(x_0) + (\hat{f}(x_0))^2\right) \\
&= \mathbb{E}\left((f(x_0))^2\right) + 2\mathbb{E}\left(\epsilon f(x_0)\right) - 2\mathbb{E}\left(f(x_0)\hat{f}(x_0)\right) + \mathbb{E}\left((\hat{f}(x_0))^2\right) + \mathrm{Var}(\epsilon) \\
&= (f(x_0))^2 - 2f(x_0)\mathbb{E}\left(\hat{f}(x_0)\right) + \mathbb{E}\left((\hat{f}(x_0))^2\right) + \mathrm{Var}(\epsilon) \ {\color{red}+} \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad {\color{red}\mathbb{E}\left(\hat{f}(x_0)\right)^2 - \mathbb{E}\left(\hat{f}(x_0)\right)^2} \\
&= \underbrace{\mathbb{E}\left((\hat{f}(x_0))^2\right) - \mathbb{E}\left(\hat{f}(x_0)\right)^2}_{\text{variance}} + \underbrace{\left(\mathbb{E}(\hat{f}(x_0)) - f(x_0)\right)^2}_{\text{bias}} + \mathrm{Var}(\epsilon)
\end{aligned}
$$

Expected **test** MSE:



$$\mathbb{E}(y_0 - \hat{f}(x_0))^2 \quad = \dots$$

$$= \mathbb{E}\big((\hat{f}(x_0))^2\big) - \mathbb{E}\big(\hat{f}(x_0)\big)^2 + \big(\mathbb{E}(\hat{f}(x_0)) - f(x_0)\big)^2 + \mathrm{Var}(\epsilon)$$

$$= \mathrm{Var}(\hat{f}(x_0)) + \big(\mathrm{Bias}(\hat{f}(x_0))\big)^2 + \mathrm{Var}(\epsilon)$$

Roughly:

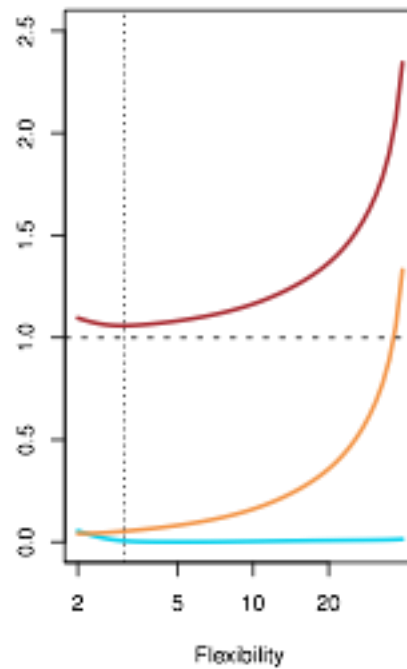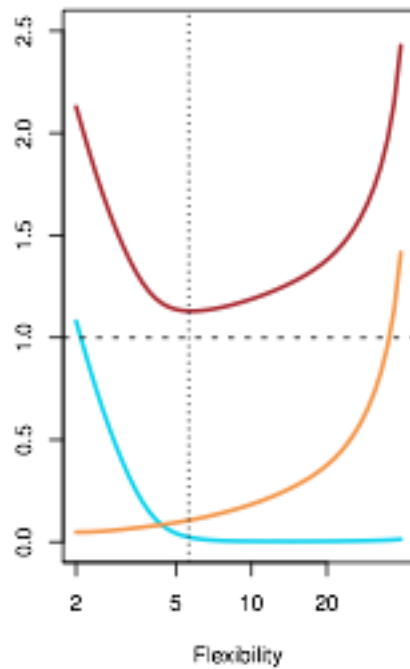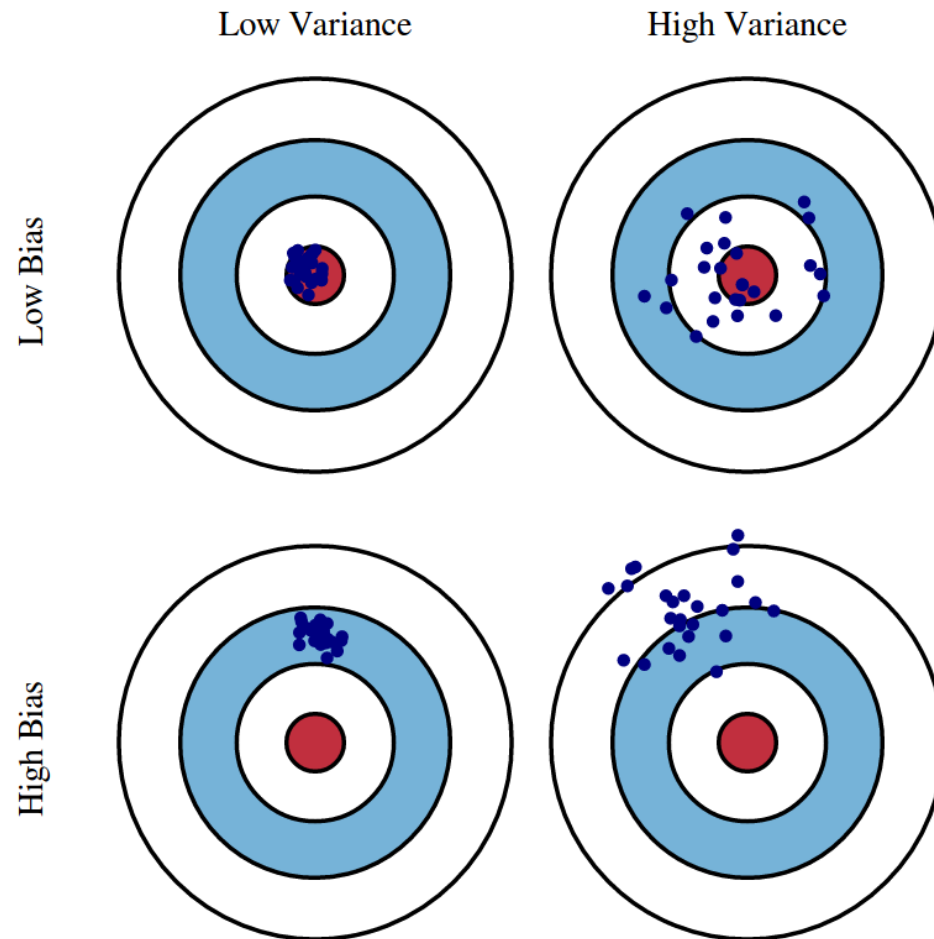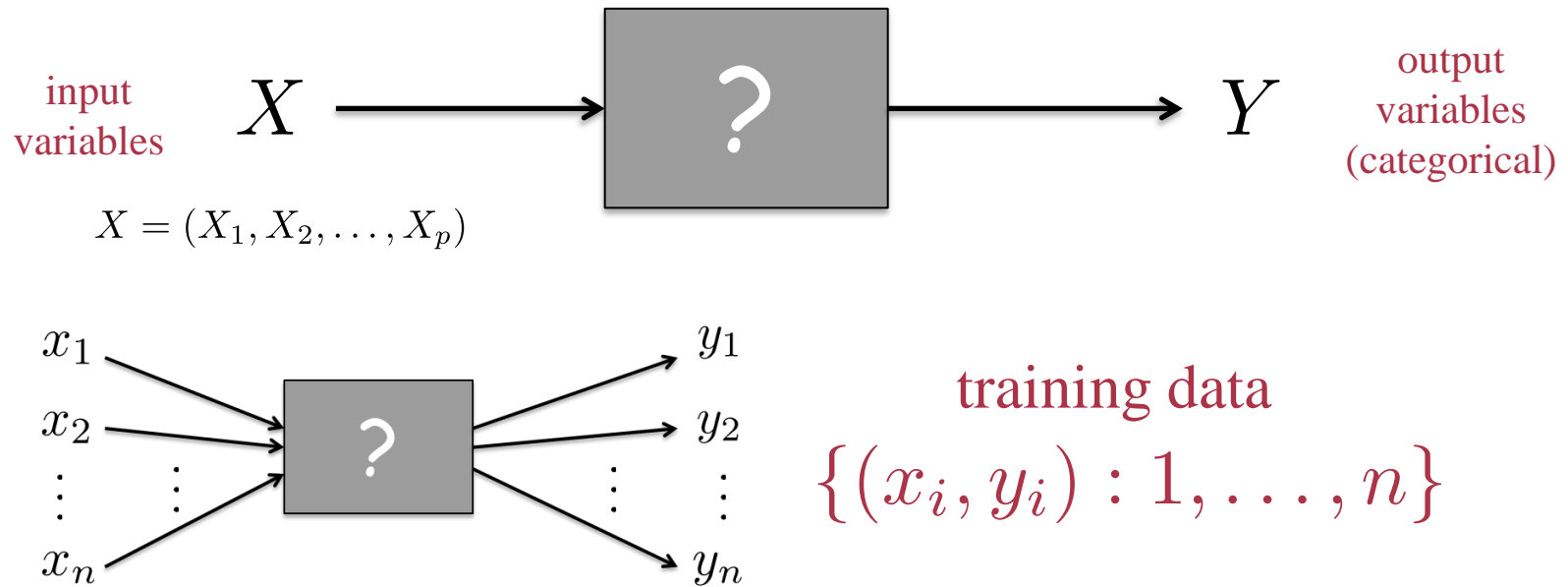| Model Complexity ⬆ | Variance ⬆ | Bias ⬇ |

$$\mathbb{E}(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + \left(\text{Bias}(\hat{f}(x_0))\right)^2 + \text{Var}(\epsilon)$$

## CLASSIFICATION

$X$  $Y$ output
variables
(categorical)

$$X = (X_1, X_2, \ldots, X_p)$$

$x_1$ $y_1$

$x_2$ ? $y_2$

$x_n$ $y_n$

training data

$$\{(x_i, y_i) : 1, \ldots, n\}$$

---

**Training** Error Rate

$$\frac{1}{n} \sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$$

---

$(x_0, y_0)$ : test data  ???

ERASMUS UNIVERSITEIT ROTTERDAM

# Confusion Matrix

| | | Predicted Class | | Total |
|---|---|---|---|---|
| | | - | + | |
| **True Class** | - | True Negative (TN) | False Positive (FP) | N |
| | + | False Negative (FN) | True Positive (TP) | P |
| | **Total** | N* | P* | |

|  | | Predicted Class | | |
| --- | --- | --- | --- | --- |
|  | | - | + | Total |
| **True Class** | - | TN | FP | N |
|  | + | FN | TP | P |
|  | **Total** | N* | P* | |

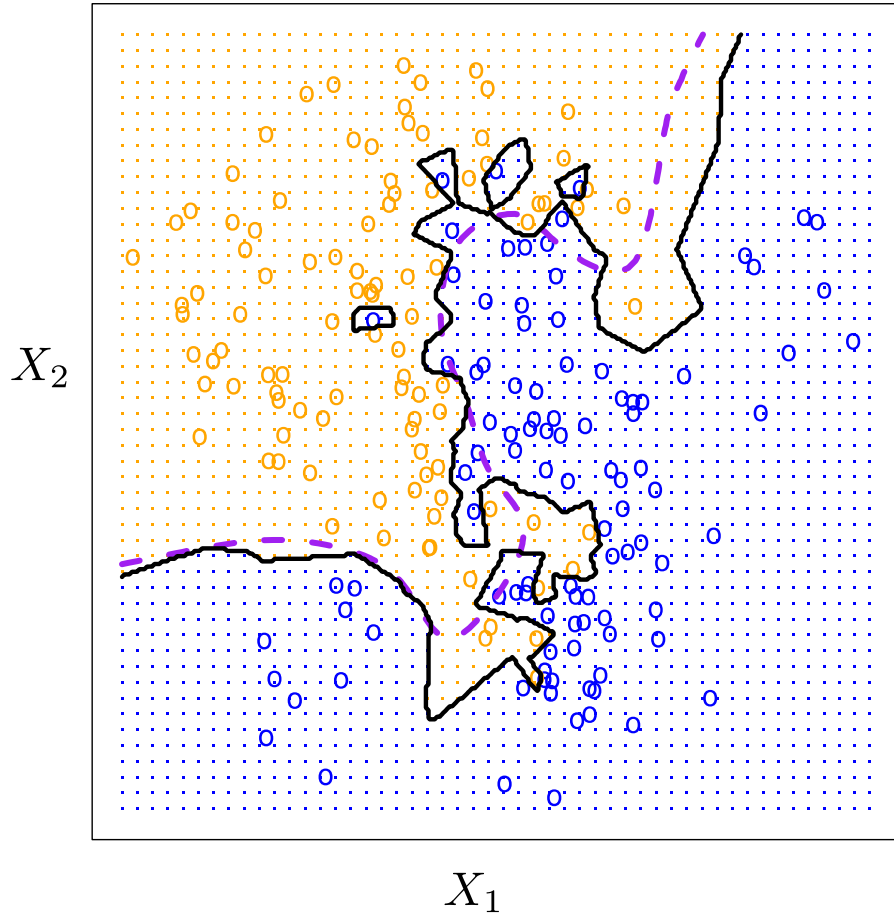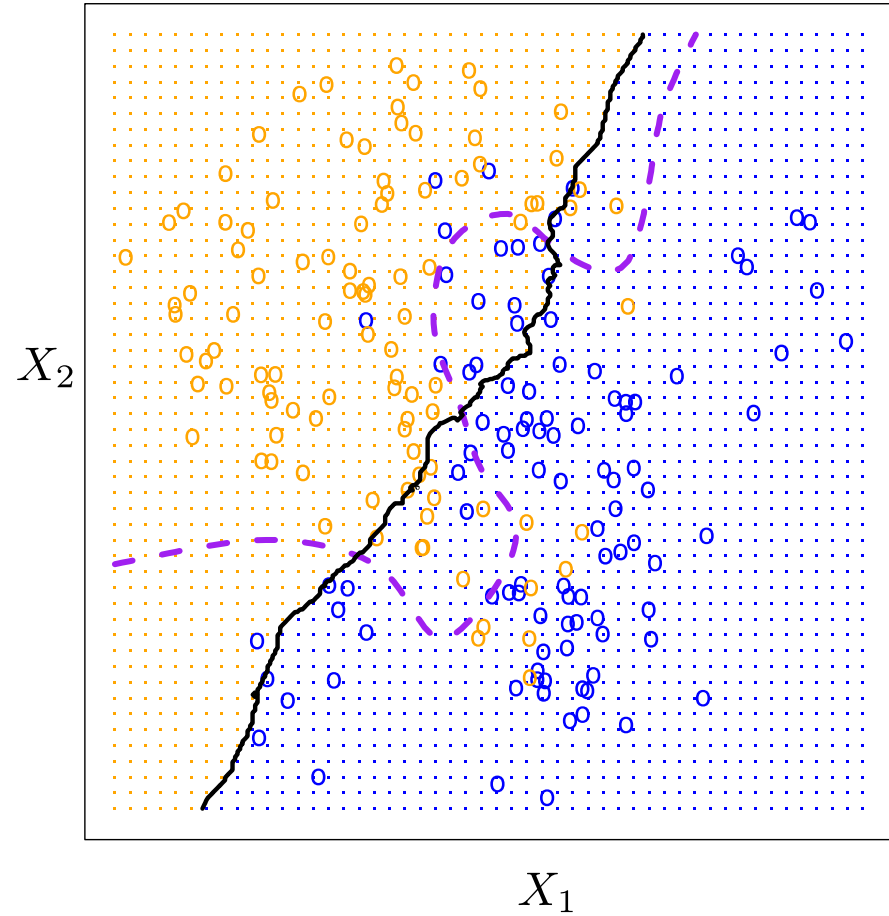| Name | Defn. | Synonyms |
| --- | --- | --- |
| **False Positive Rate** | FP/N | type I error, (1-specificity) |
| **True Positive Rate** | TP/P | (1-type II error), power, sensitivity, recall |
| **Positive Pred. Value** | TP/P* | precision, (1-false discovery proportion) |
| **Negative Pred. Value** | TN/N* | |

# $K$ – Nearest Neighbors



$$\mathbb{P}(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j)$$

majority voting

set of $K$
closest points to $x_0$

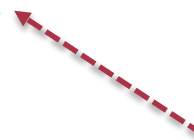$K = 1$

$K = 100$

**Voronoi Tessellation**

# Regression

$X, Y$ random: $f(X) =$?

$$\mathbb{E}((Y - f(X))^2) = \mathbb{E}(\mathbb{E}((Y - f(X))^2 \mid X))$$

$$f(x) = \arg\min_u \underbrace{\mathbb{E}((Y - u)^2 \mid X = x)}_{h(u)}$$

$$\frac{\partial h(u)}{\partial u} = 2u - 2\mathbb{E}(Y|X = x) = 0 \implies f(x) = \mathbb{E}(Y|X = x)$$

$$\boxed{f(x_0) = \mathbb{E}(Y|X = x_0)}$$

$$\hat{f}(x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} y_i$$

*K*-NN approximation to regression

ERASMUS UNIVERSITEIT ROTTERDAM

# Classification

$X, Y$ random: $\hat{Y}(X) = ?$

$$\mathbb{E}(I(Y \neq \hat{Y}(X))) = \mathbb{E}(\mathbb{E}(I(Y \neq \hat{Y}(X)) \mid X))$$

$$\hat{Y}(x) = \arg\min_{j} \; \mathbb{E}(I(Y \neq j) \mid X = x)$$

$$= \arg\min_{j} \; \sum_{k} I(k \neq j)\mathbb{P}(Y = k | X = x)$$

$$= \arg\min_{j} \; (1 - \mathbb{P}(Y = j | X = x))$$

$$= \arg\max_{j} \; \mathbb{P}(Y = j | X = x)$$

# Bayes Classifier

$$\hat{y}_0 = \arg\max_{j} \ \mathbb{P}(Y = j | X = x_0)$$

## Bayes Error Rate

$$\mathbb{E}\Big(1 - \max_{j} \ \mathbb{P}(Y = j | X)\Big) = 1 - \mathbb{E}\Big(\max_{j} \ \mathbb{P}(Y = j | X)\Big)$$