# Machine learning primer
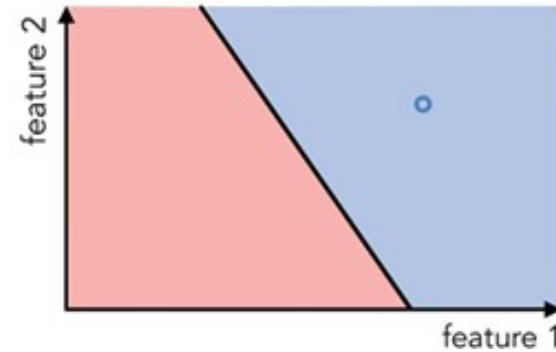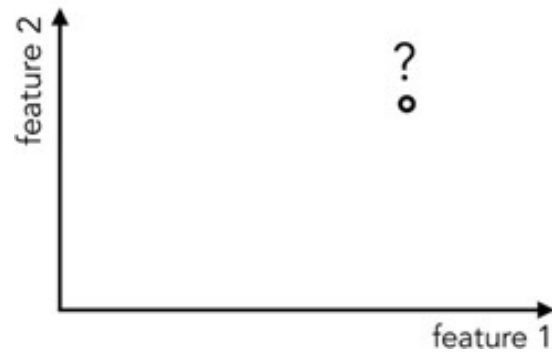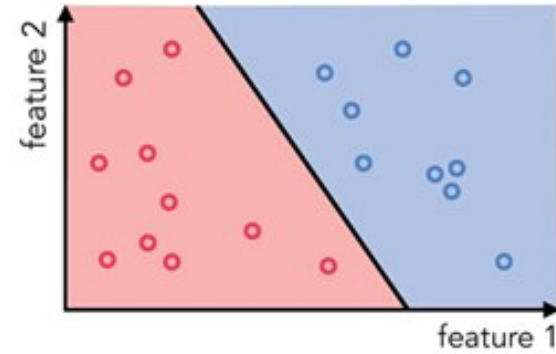
JOANNA BYSZUK & JEREMI OCHAB
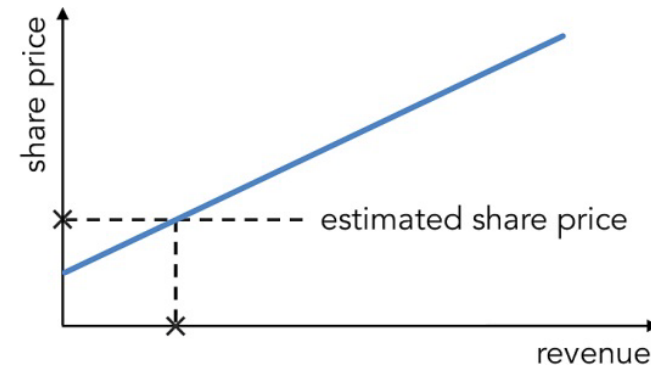
DHSI 2024, "DIY COMPUTATIONAL TEXT ANALYSIS WITH R"
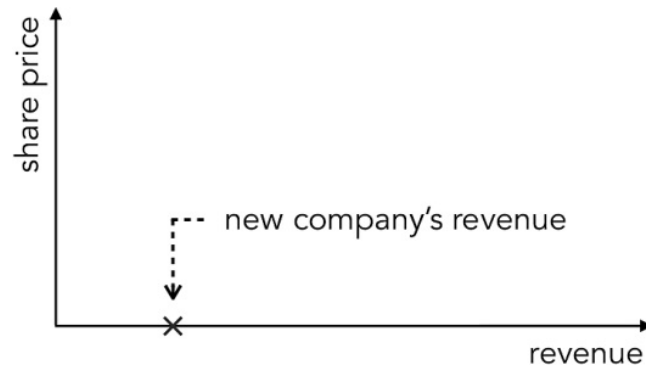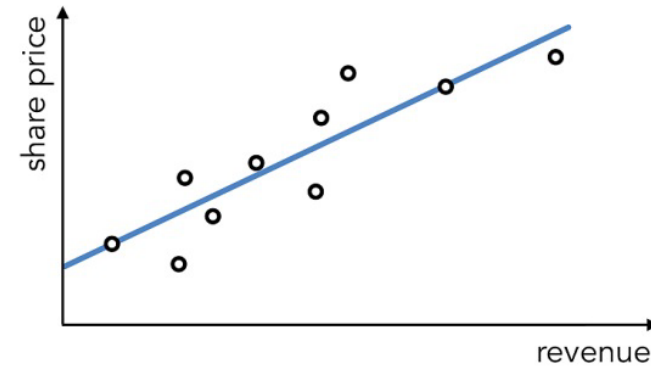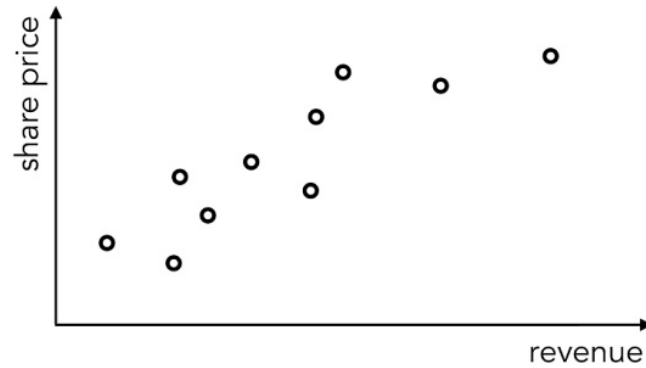
# Task examples: classification

# Task examples: regression

# Terms and definitions

**Example**: item, instance of the data used.

**Features**: attributes associated to an item, often represented as a vector (e.g., word counts).

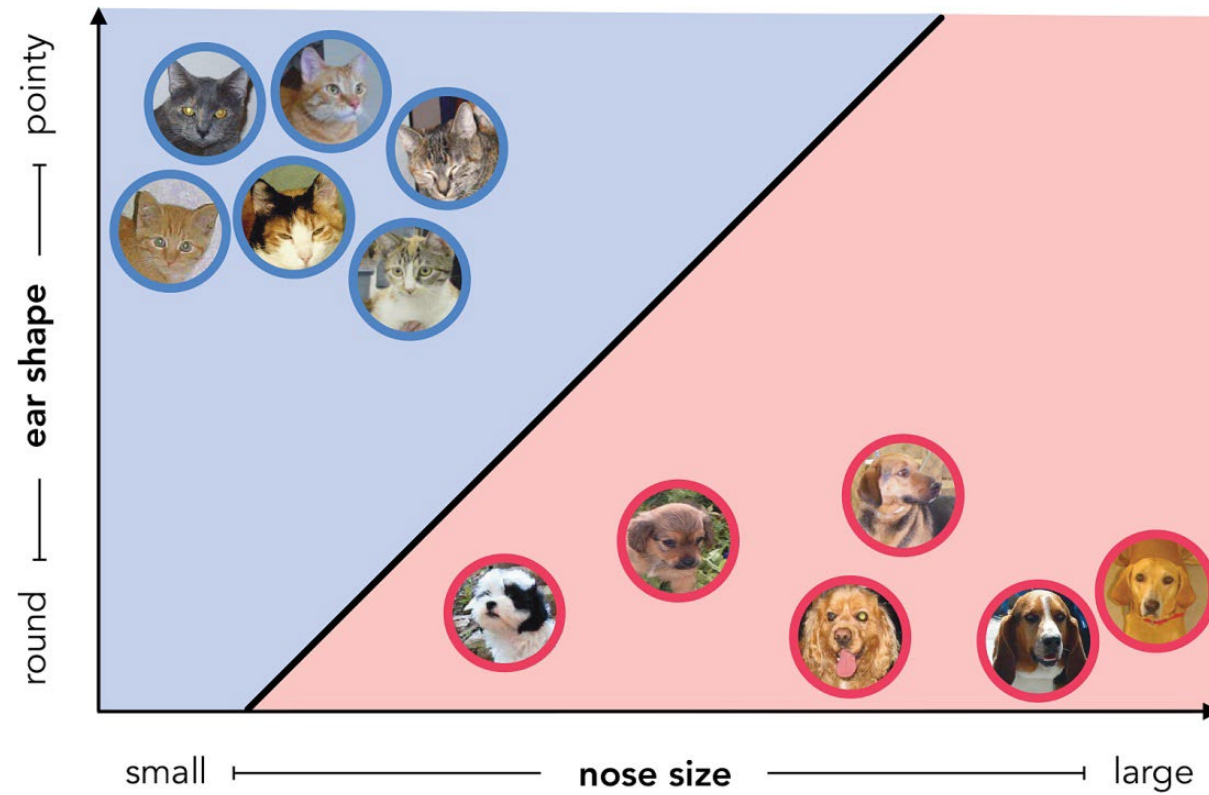**Labels**: category (classification) or real value (regression) associated to an item.

**Data**:

❖training data (typically labeled).

❖test data (labeled but labels not seen).
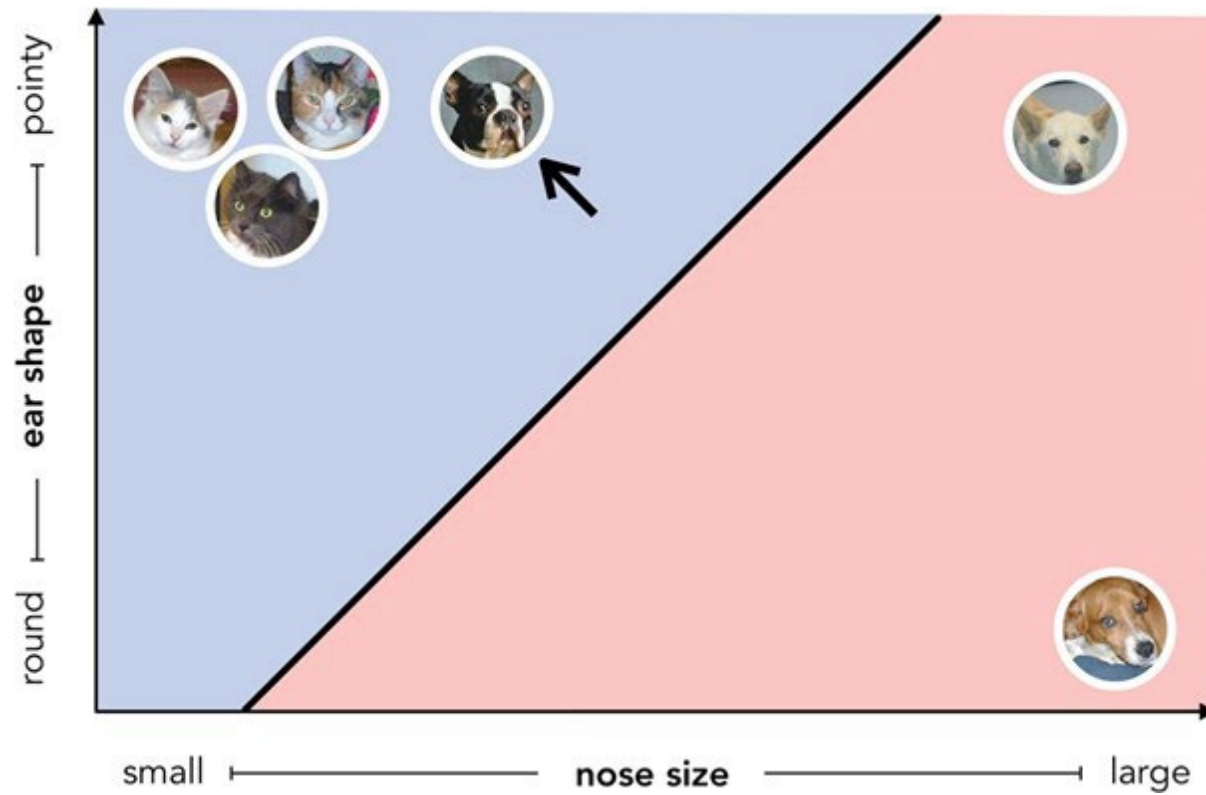
❖validation data
(labeled, for tuning parameters).

# Learning stages

# Learning stages

# Learning stages

# Learning stages



Data collection → Feature design → Model training → **Model selection** → **Model evaluation**

Training set

Validation set

**Test set**

$A(\Theta)$

$A(\Theta_0)$

previous knowledge

algorithm

training/empirical error

validation/off-training set error

test (~generalization) error

# Bias-variance tradeoff

# Bias-variance tradeoff
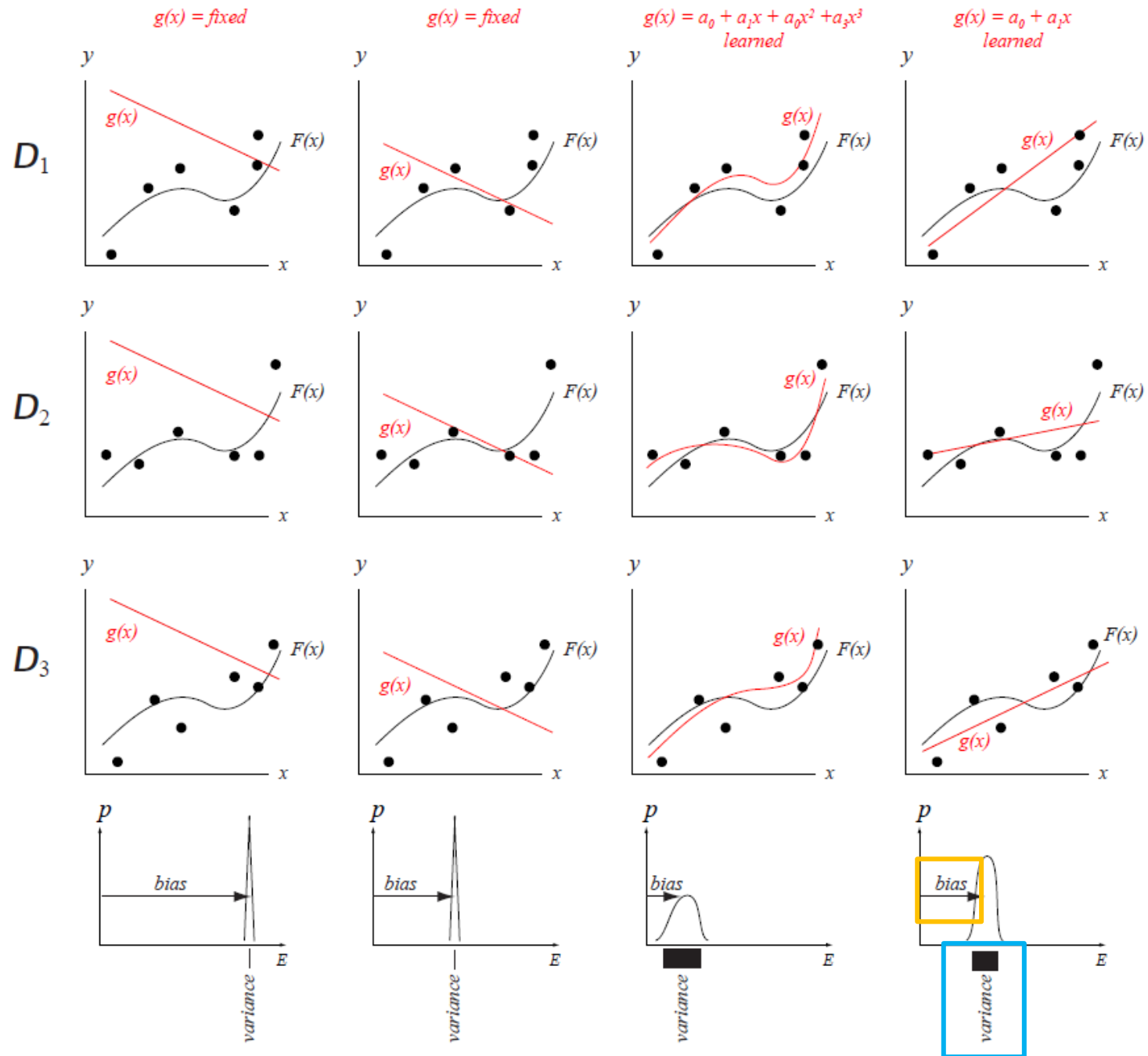
Najłatwiej na przykładzie dopasowania krzywej regresji:

- $F(\mathbf{x})$ – a true (but unknown) function with continuous valued output with noise

- $D$ – set of $n$ training samples generated by $F(\mathbf{x})$
- $g(\mathbf{x}; D)$ – the estimated regressions function F (depends on the set $D$!)

- Estimator effectiveness: average (over all sets $D$ of size $n$) mean-square deviation:

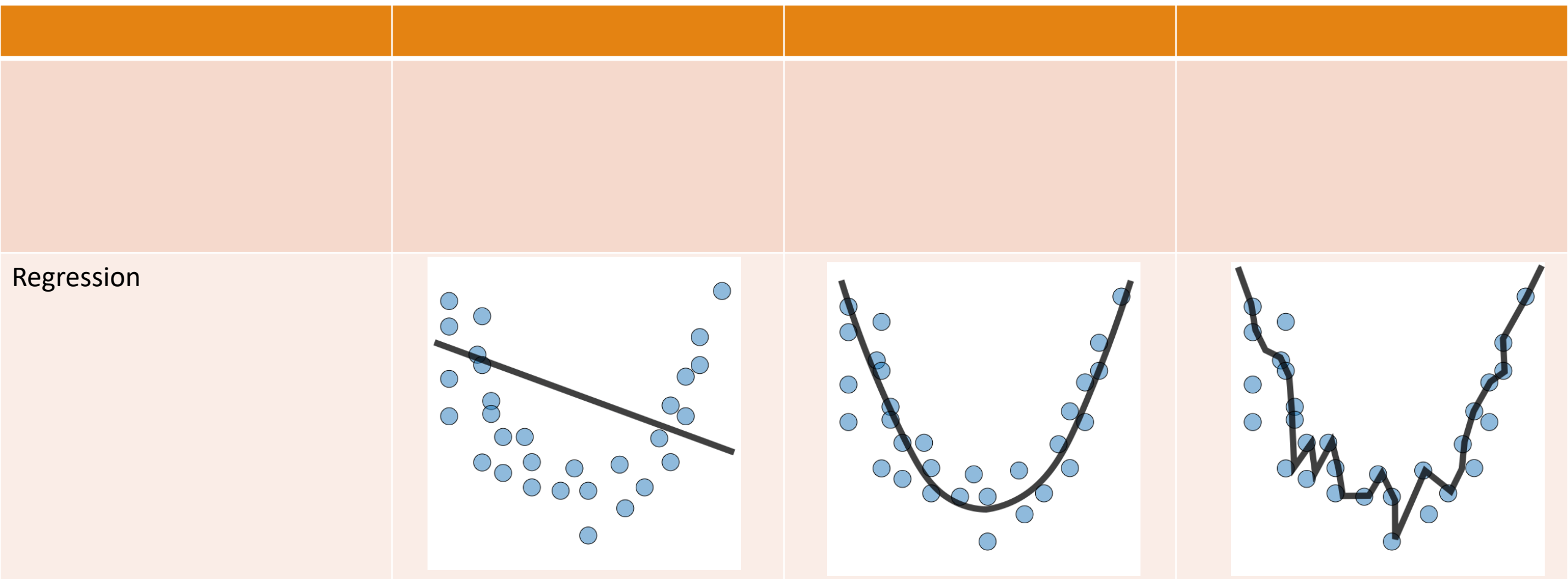○ bias – the difference between the true (but unknown) value and our expectations [=estimation accuracy]

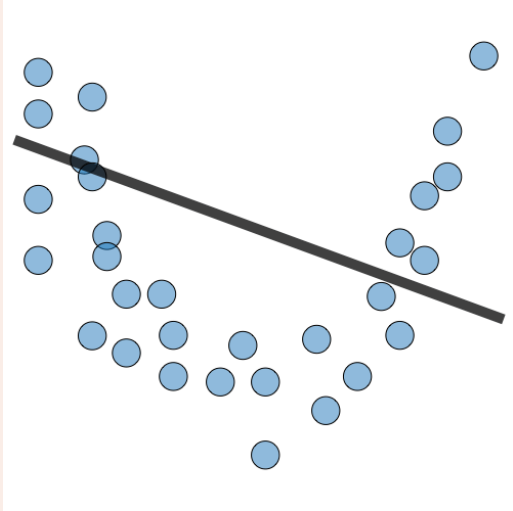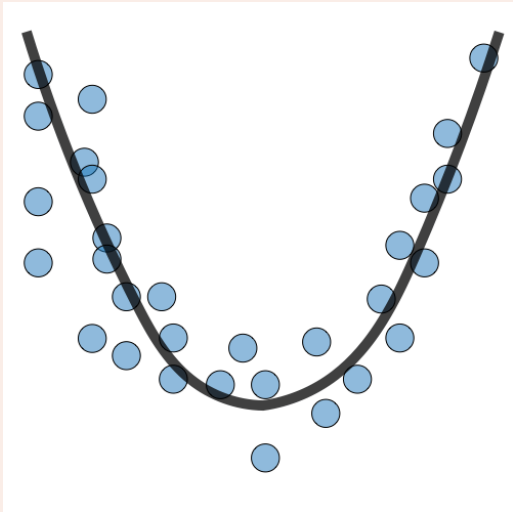○ variance – the instability of the estimate due to the variability of the training set

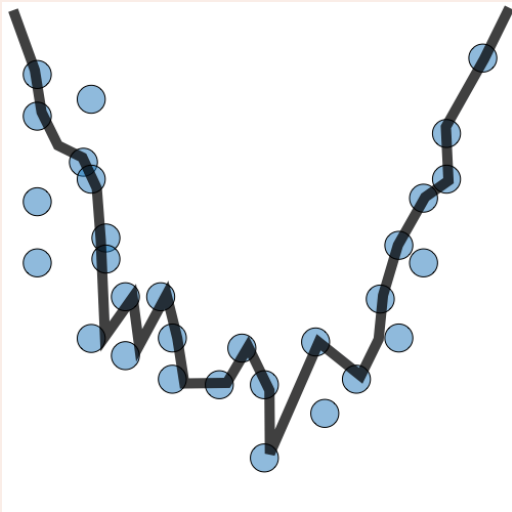$$\mathcal{E}_D\left[(g(\mathbf{x};\ \mathcal{D}) - F(\mathbf{x}))^2\right]$$
$$= \underbrace{(\mathcal{E}_D[g(\mathbf{x};\ \mathcal{D}) - F(\mathbf{x})])^2}_{bias^2} + \underbrace{\mathcal{E}_D\left[(g(\mathbf{x};\ \mathcal{D}) - \mathcal{E}_D[g(\mathbf{x};\ \mathcal{D})])^2\right]}_{variance}.$$

Duda, Hart, Stork. *Pattern Classification, 2nd ed*. Wiley (2000).

E – mean squared error

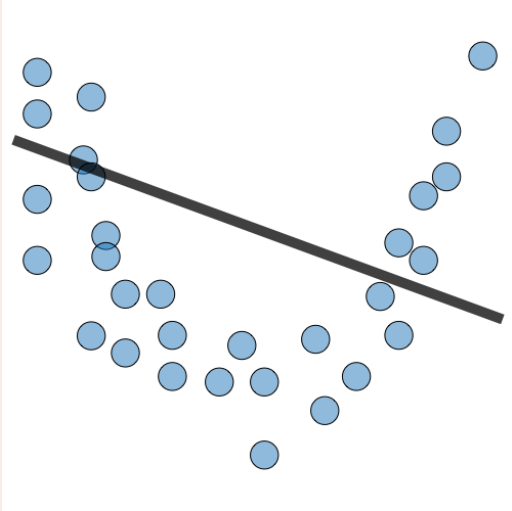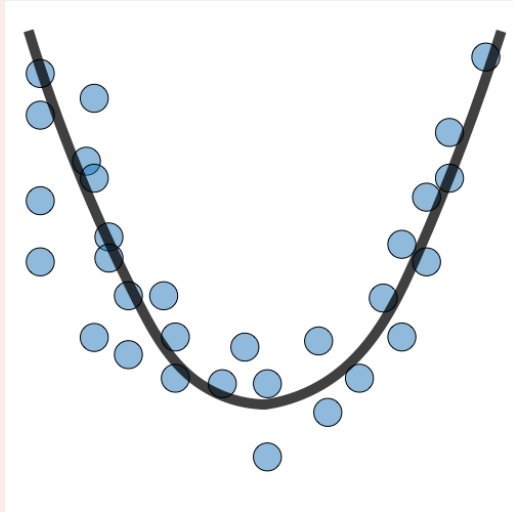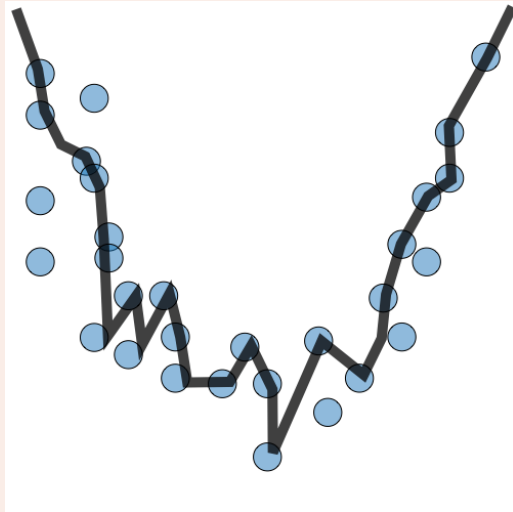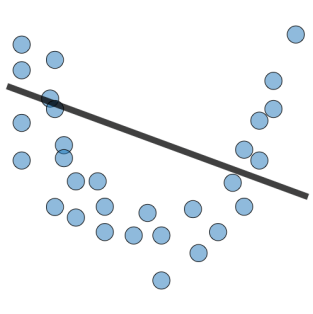p – probability that we will randomly get *D* with the error E

# Bias-variance tradeoff

| | | | |
|---|---|---|---|
| | | | |
| Regression |  |  |  |

# Bias-variance tradeoff

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| | | | |
| Regression |  |  |  |

Afshine Amidi, Shervine Amidi. *Machine Learning cheatsheets for Stanford's CS 229*. (2018).

# Bias-variance tradeoff

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| Symptoms | <ul><li>High training error</li><li>Training error close to test error</li><li>High bias</li></ul> | <ul><li>Training error slightly lower than test error</li></ul> | <ul><li>Low training error</li><li>Training error much lower than test error</li><li>High variance</li></ul> |
| Regression |  |  |  |

| | Underfitting | W sam raz | Overfitting |
|---|---|---|---|
| Symptoms | ▪ High training error<br>▪ Training error close to test error<br>▪ High bias | ▪ Training error slightly lower than test error | ▪ Low training error<br>▪ Training error much lower than test error<br>▪ High variance |
| Regression |  |  |  |
| Classification |  |  |  |
| Deep learning |  |  |  |
| Remedies? | ▪ complexify model<br>▪ Add more features<br>▪ Train longer | | ▪ Regularise<br>▪ Get more data |

# Generalisation
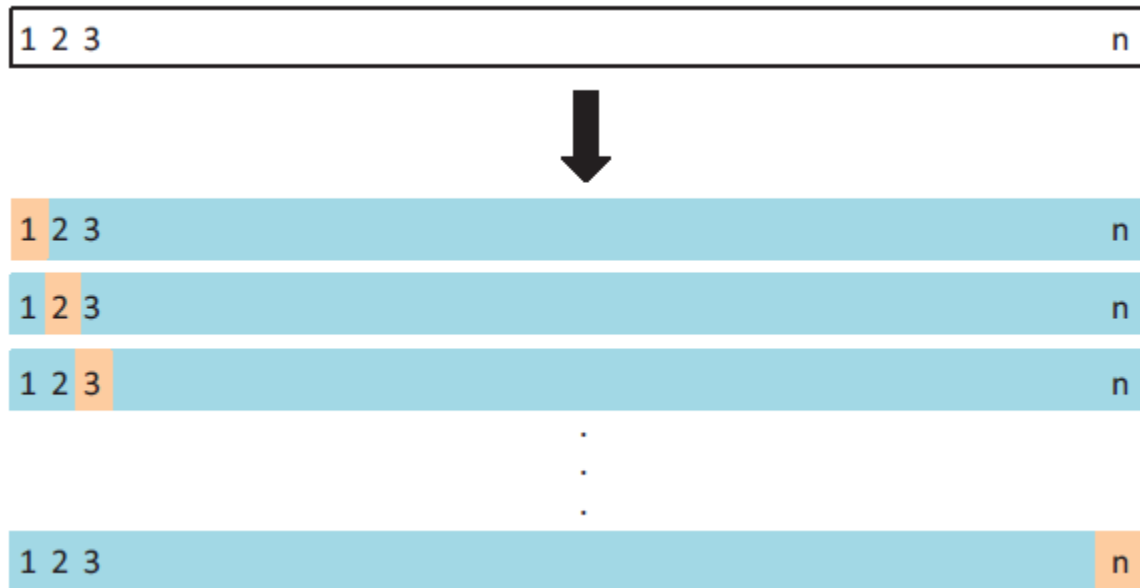
❖the best hypothesis on the sample may not be the best overall.

❖generalization is not memorization.

❖complex rules (very complex separation surfaces) can be poor predictors.

❖trade-off: complexity of hypothesis set vs sample size (underfitting/overfitting).
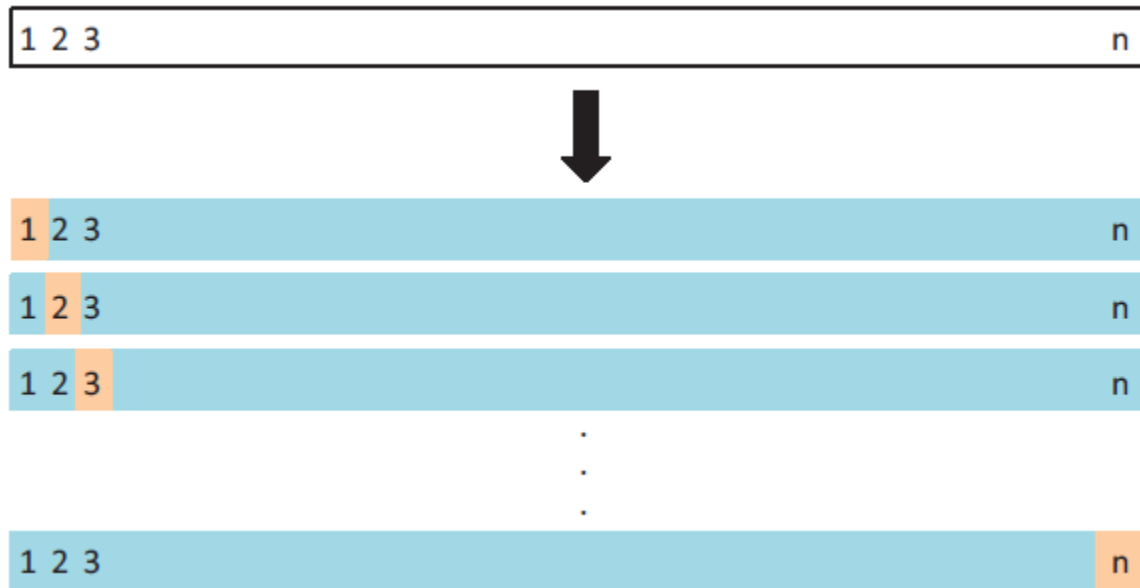
# Cross-validation

James, Witten, Hastie, Tibshirani, *An Introduction to Statistical Learning,* Springer (2013).
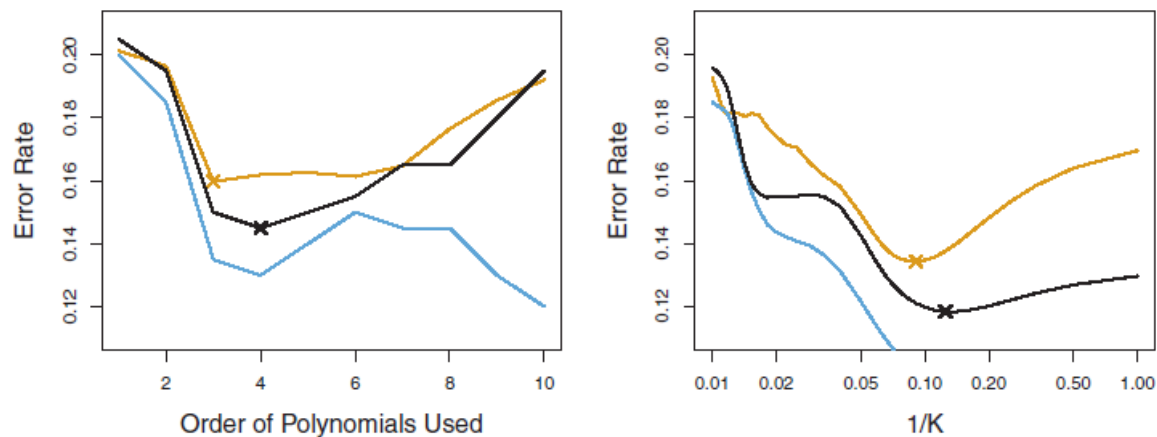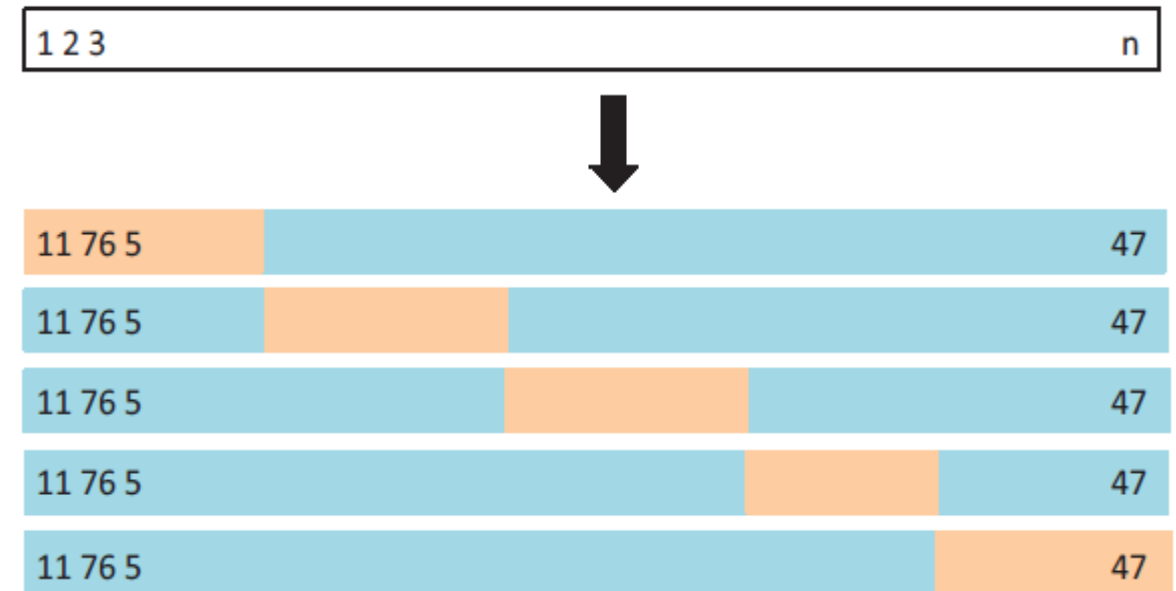
# Cross-validation

# Cross-validation

FIGURE 5.8. *Test error (brown), training error (blue), and 10-fold CV error (black) on the two-dimensional classification data displayed in Figure 5.7. Left: Logistic regression using polynomial functions of the predictors. The order of the polynomials used is displayed on the x-axis. Right: The KNN classifier with different values of K, the number of neighbors used in the KNN classifier.*

James, Witten, Hastie, Tibshirani, *An Introduction to Statistical Learning,* Springer (2013).

# Cross-validation

- Stratified CV                                            stratified

- Group CV

[https://scikit-learn.org/stable/modules/cross_validation.html#k-fold](https://scikit-learn.org/stable/modules/cross_validation.html#k-fold)