# model overfitting

underfitting = model is too simple, both training and test errors are large

overfitting = model is too complex, training error is small, but test error is large

• increasing the size of training data reduces the differences between training and testing errors

reasons for overfitting = not enough training data or high model complexity

multiple comparison procedure =

| Day 1 | Up |
| Day 2 | Down |
| Day 3 | Down |
| Day 4 | Up |
| Day 5 | Down |
| Day 6 | Down |
| Day 7 | Up |
| Day 8 | Up |
| Day 9 | Up |
| Day 10 | Down |

$$P(\text{number of correct guesses} \geq 8) = \frac{\binom{10}{8} + \binom{10}{9} + \binom{10}{10}}{2^{10}} = 0.0547$$

all 50 not guess more than 8

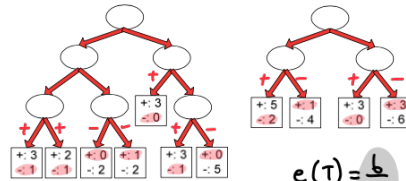$$P(\text{among 50 analyst, at least one makes at least 8 correct guesses}) = 1 - (1 - 0.0547)^{50}$$
$$= 0.9399$$

generalization error (model) = training error (model, training data) + $\alpha \times$ complexity (model)

complexity of decision trees → # of leaf nodes

## pessimistic error estimate =

$$err_{gen}(T) = err(T) + \Omega \times \frac{k}{N_{train}}$$

− err(T): error rate on all training records
− $\Omega$: trade-off hyper-parameter (similar to $\alpha$)
  ◆ Relative cost of adding a leaf node
− k: number of leaf nodes
− $N_{train}$: total number of training records



| +: 3 | +: 2 | +: 0 | +: 1 | | +: 3 | +: 0 |
| -: 1 | -: 1 | -: 2 | -: 2 | | -: 1 | -: 5 |

+: 3
-: 0

$$e(T) = \frac{4}{24}$$

$$e_{gen}(T) = \frac{4}{24} + 1 \cdot \frac{7}{24} = \frac{11}{24}$$

| +: 5 | +: 1 | | +: 3 | +: 3 |
| -: 2 | -: 4 | | -: 0 | -: 6 |

$$e(T) = \frac{6}{24}$$

$$e_{gen}(T) = \frac{6}{24} + 1 \cdot \frac{4}{24} = \frac{10}{24}$$

→ optimistic errors

optimistic error estimate = $err_{gen}(T) = err(T)$ using only training error = resubstitution estimate

minimum description length (MDL) = $\underbrace{Cost(Model, Data)}_{\text{# of bits needed for encoding}} = \underbrace{Cost(Data | Model)}_{\text{encoding of misclassification err}} + \alpha \underbrace{Cost(Model)}_{\text{encoding of nodes + split condition}}$