# Introduction to Machine Learning

**Multi-Class Classification
and
Bayesian optimization**

# Multiclass classification

- More than two classes
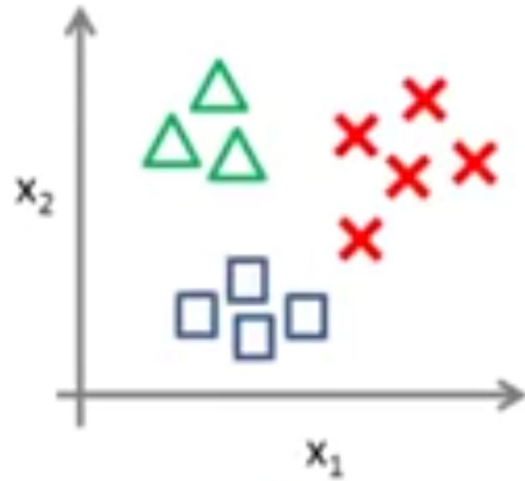
# Confusion matrix

Predict →

| | Enter | Leave | Cook | Sleep | Meds | Eat | Groom | Bathe | Bed-T | Relax |
|---|---|---|---|---|---|---|---|---|---|---|
| Enter | 1673 | 27 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Leave | 9 | 1979 | 5 | 2 | 0 | 1 | 0 | 0 | 0 | 0 |
| Cook | 59 | 58 | 51238 | 39 | 199 | 137 | 28 | 2 | 0 | 0 |
| Sleep | 21 | 29 | 5 | 30795 | 4 | 86 | 14 | 0 | 51 | 0 |
| Meds | 11 | 2 | 200 | 0 | 3105 | 1 | 0 | 0 | 0 | 0 |
| Eat | 3 | 3 | 6 | 94 | 1 | 14278 | 5 | 0 | 0 | 0 |
| Groom | 0 | 11 | 4 | 1 | 1 | 3 | 21833 | 33 | 41 | 0 |
| Bathe | 0 | 0 | 0 | 1 | 0 | 0 | 59 | 592 | 5 | 0 |
| Bed-T | 0 | 0 | 0 | 18 | 0 | 0 | 15 | 2 | 501 | 0 |
| Relax | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 3 |

Actual →

# One-versus-all (one-vs-rest)



Class 1: △
Class 2: □
Class 3: ✕

**Algorithm 13** $\textsc{OneVersusAllTrain}(\mathbf{D}^{multiclass}, \textsc{BinaryTrain})$

1:  **for** $i = 1$ **to** $K$ **do**
2:      $\mathbf{D}^{bin} \leftarrow$ relabel $\mathbf{D}^{multiclass}$ so class $i$ is positive and $\neg i$ is negative
3:      $f_i \leftarrow \textsc{BinaryTrain}(\mathbf{D}^{bin})$
4:  **end for**
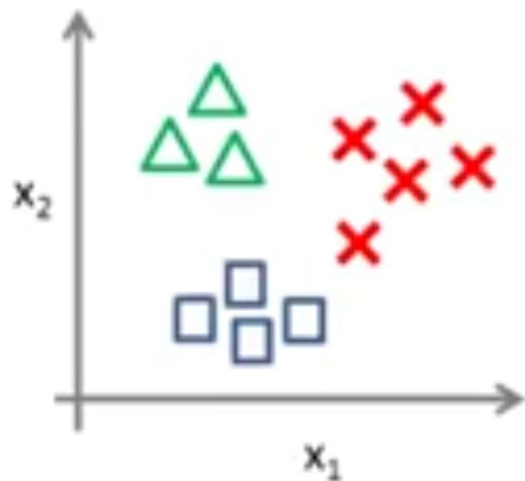5:  **return** $f_1, \dots, f_K$

---

**Algorithm 14** $\textsc{OneVersusAllTest}(f_1, \dots, f_K, \hat{x})$

1:  $score \leftarrow \langle 0, 0, \dots, 0 \rangle$                    // initialize $K$-many scores to zero
2:  **for** $i = 1$ **to** $K$ **do**
3:      $y \leftarrow f_i(\hat{x})$
4:      $score_i \leftarrow score_i + y$
5:  **end for**
6:  **return** $\operatorname{argmax}_k score_k$

# All-versus-all (one-versus-one)



Class 1: △
Class 2: □
Class 3: ✖

**Algorithm 15** AllVersusAllTrain($\mathbf{D}^{multiclass}$, BinaryTrain)

1: $f_{ij} \leftarrow \emptyset, \forall 1 \leq i < j \leq K$
2: **for** $i = 1$ **to** $K\text{-}1$ **do**
3:      $\mathbf{D}^{pos} \leftarrow$ all $x \in \mathbf{D}^{multiclass}$ labeled $i$
4:      **for** $j = i+1$ **to** $K$ **do**
5:          $\mathbf{D}^{neg} \leftarrow$ all $x \in \mathbf{D}^{multiclass}$ labeled $j$
6:          $\mathbf{D}^{bin} \leftarrow \{(x, +1) : x \in \mathbf{D}^{pos}\} \cup \{(x, -1) : x \in \mathbf{D}^{neg}\}$
7:          $f_{ij} \leftarrow$ BinaryTrain($\mathbf{D}^{bin}$)
8:      **end for**
9: **end for**
10: **return** all $f_{ij}$s

---

**Algorithm 16** AllVersusAllTest(all $f_{ij}$, $\hat{x}$)

1: $score \leftarrow \langle 0, 0, \ldots, 0 \rangle$                        // initialize $K$-many scores to zero
2: **for** $i = 1$ **to** $K\text{-}1$ **do**
3:      **for** $j = i+1$ **to** $K$ **do**
4:          $y \leftarrow f_{ij}(\hat{x})$
5:          $score_i \leftarrow score_i + y$
6:          $score_j \leftarrow score_j - y$
7:      **end for**
8: **end for**
9: **return** $\text{argmax}_k \ score_k$

# Binary tree of classifiers

# Overrun by hyperparameters

- Manual
- Grid search
- Random search

# Bayesian optimization to the rescue?

- Uses Bayes Theorem to direct the search

$$\text{Probability(event)} = P(\text{event}) = \frac{\text{\#instances of the event}}{\text{total \#instances}}$$
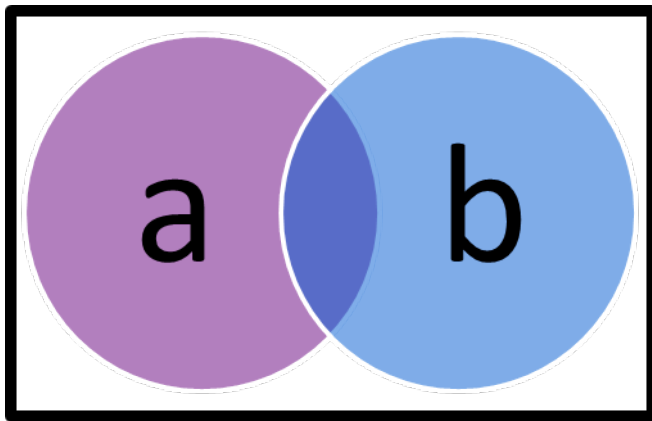
# Roll two dice

# Sources of probabilities

- Frequency
- Consider the probability that the sun will still exist tomorrow.
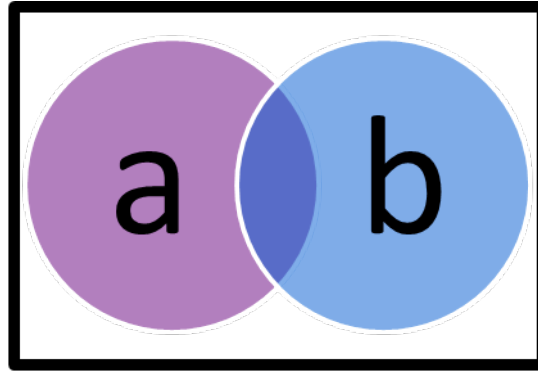
# Axioms of probability

- $0 \leq P(\text{Event}) \leq 1$

- Disjunction, $P(a \text{ or } b) = P(a) + P(b) - P(a \text{ and } b)$

# Conditional probability and conjunction

- P(a|b) = P(a and b) / P(b)

# Conditional probability and conjunction

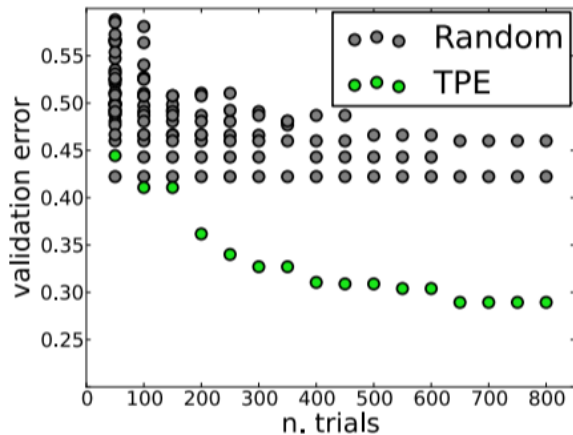- P(a and b) = P(a) $\times$ P(b|a)
- P(a and b) = P(b) $\times$ P(a|b)

- If a and b are independent events
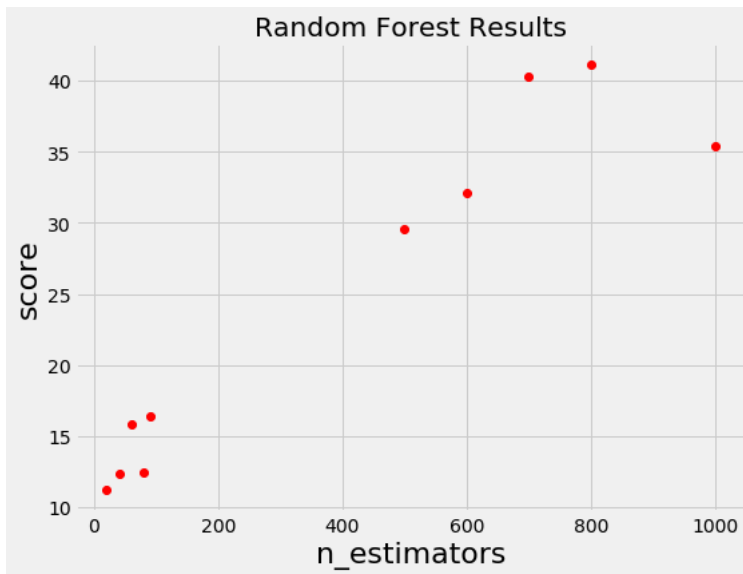  - P(a and b) = P(a) $\times$ P(b)

# Bayes' rule

# Bayesian optimization to the rescue?
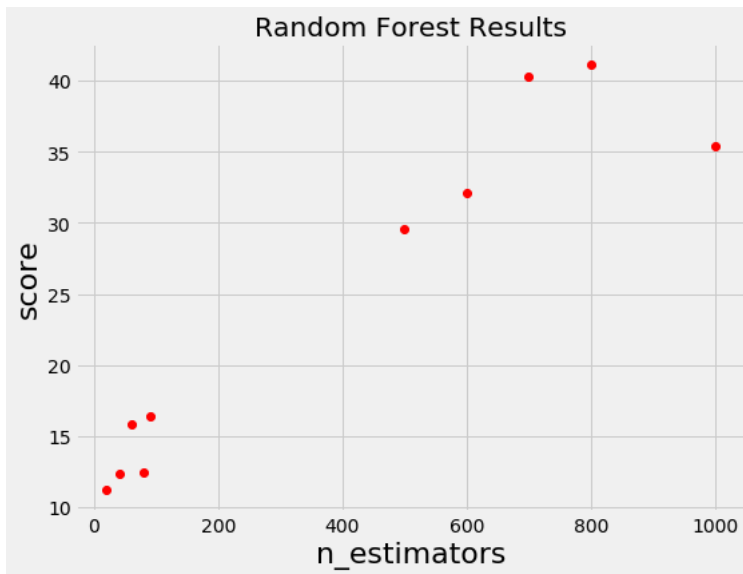
- Optimization method to solve $\arg \min_{x \in X} f(x)$

# Bayesian optimization

- Build probability model of objective function
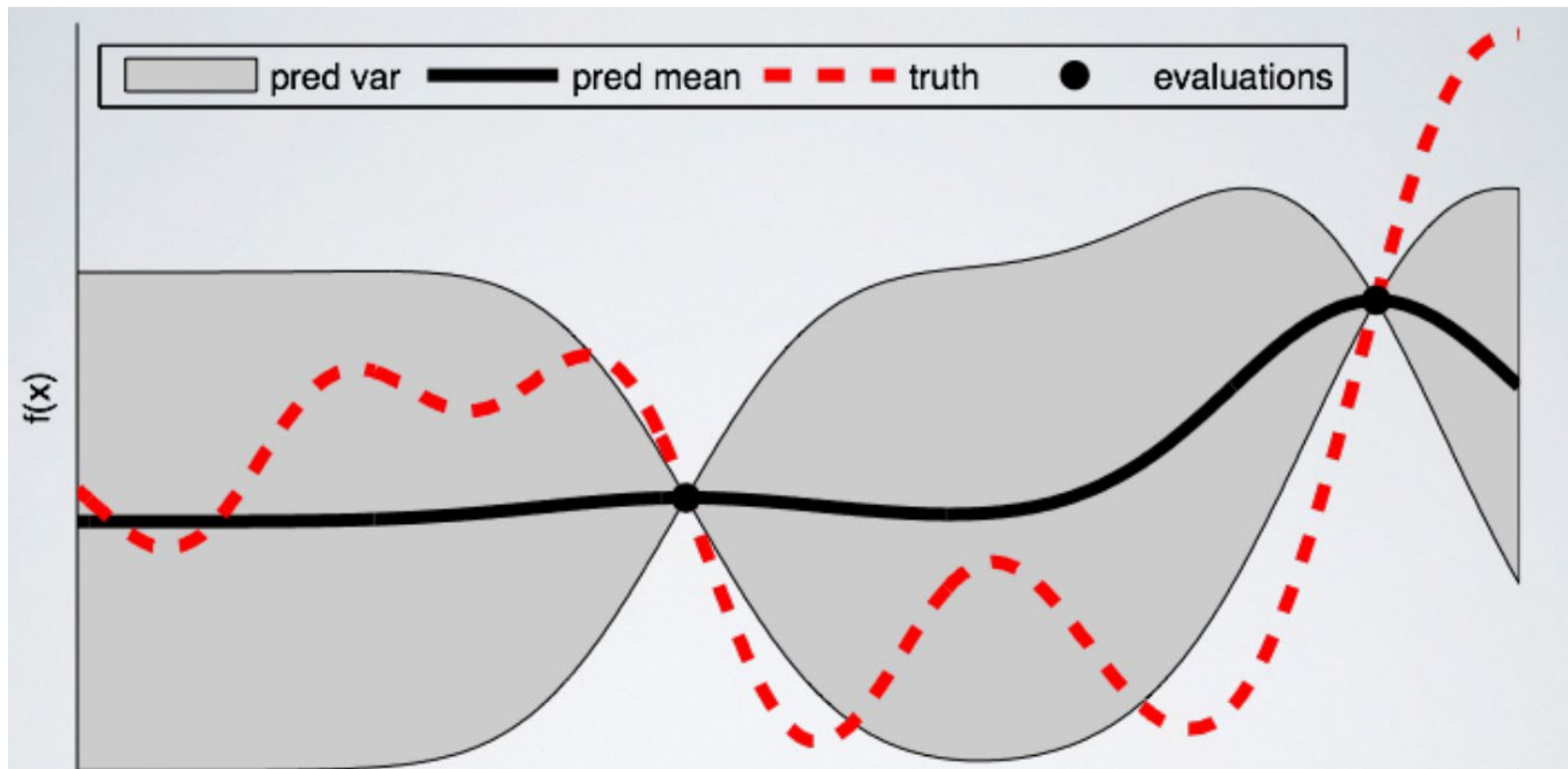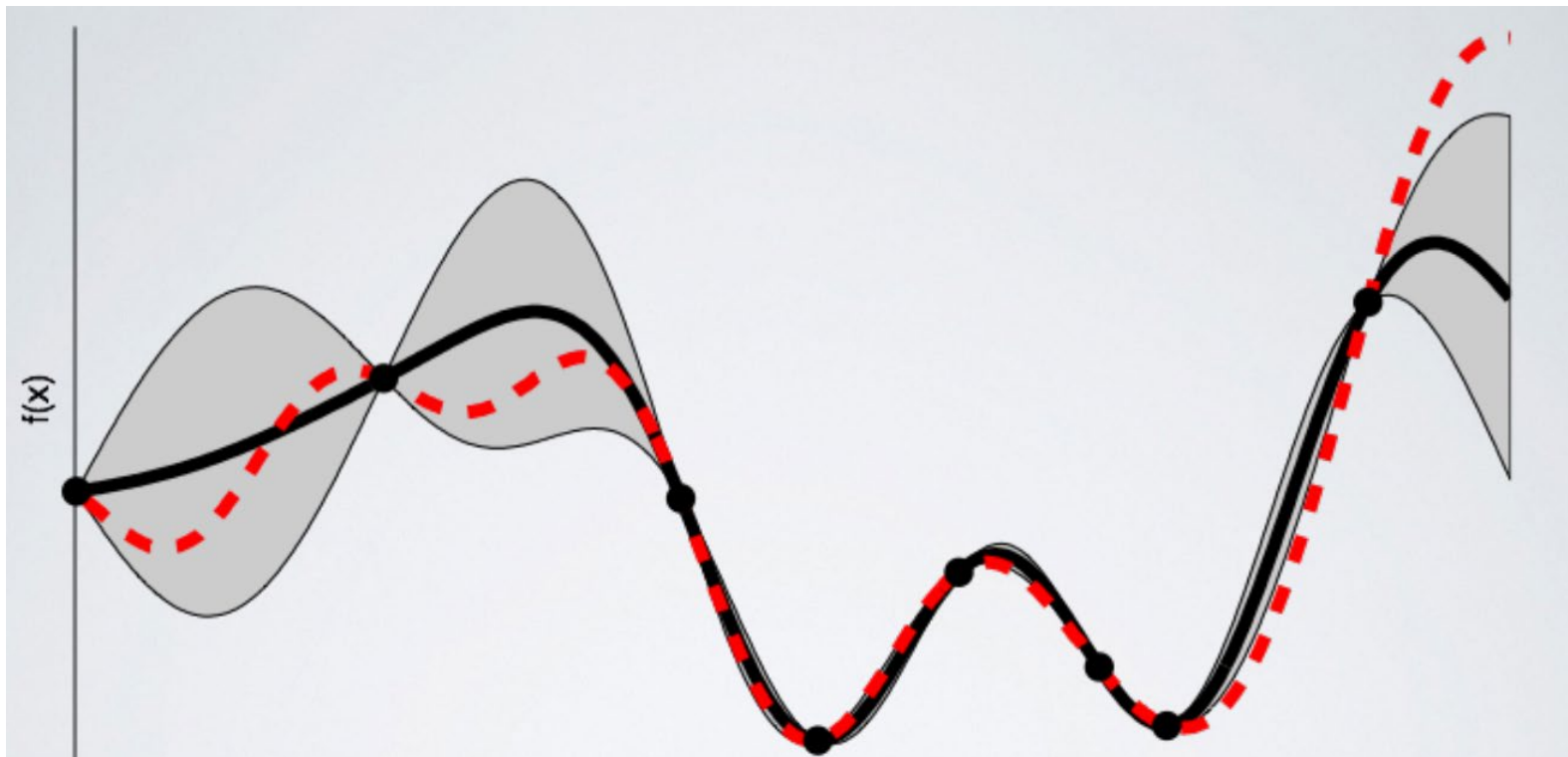- Use model to select hyperparameters to evaluate



Random Forest Results

# Bayesian optimization

- P(score | hyperparameters)

# Surrogate model

# Surrogate model

# Simple 1D example

# Compare Bayesian optimizer with random search
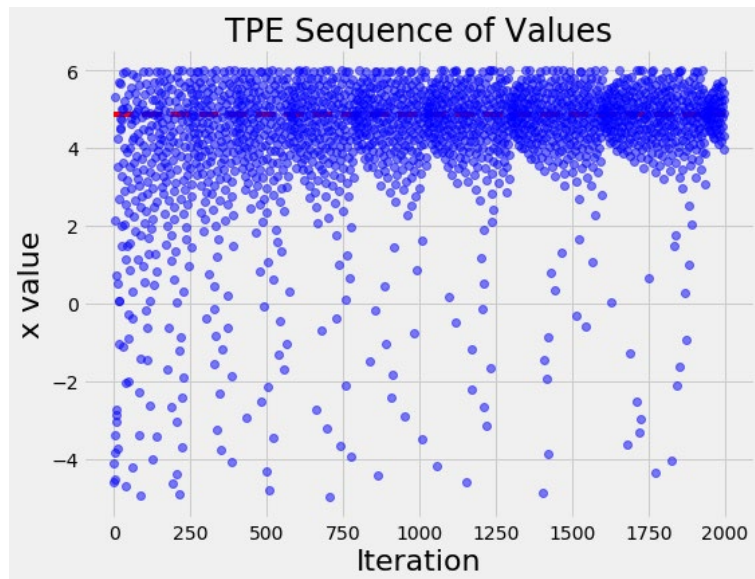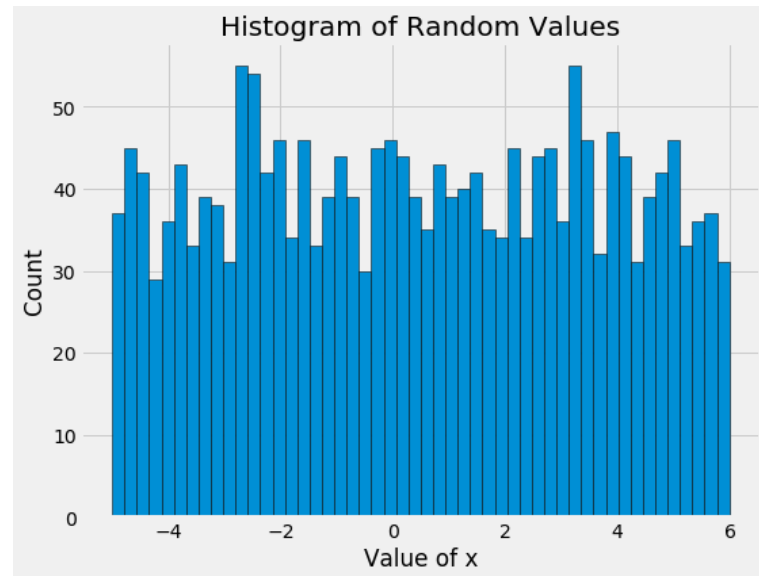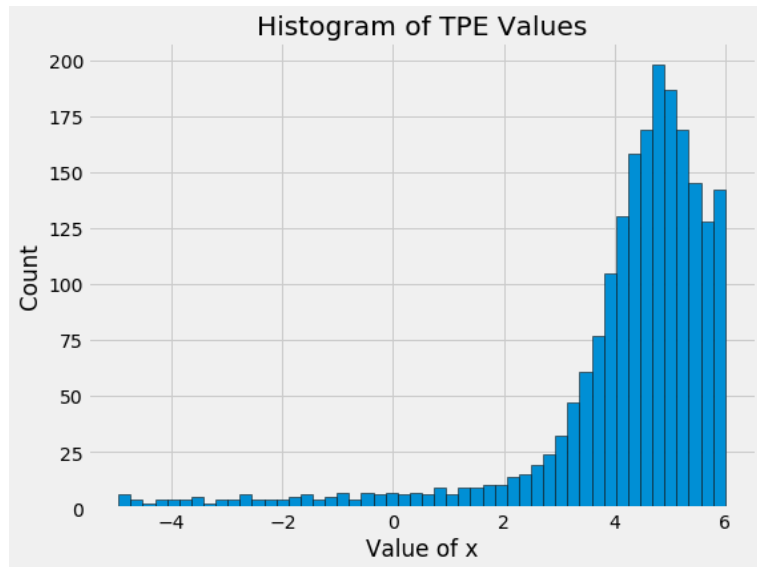
# Compare Bayesian optimizer with random search

Let's try it out