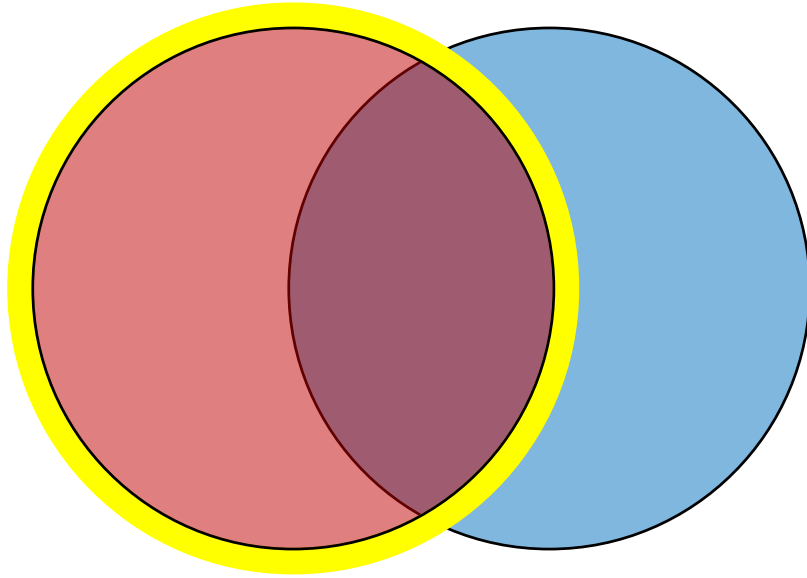# Chapter 6

# Naïve Bayes
# Grid Search & Pipelines

# Naïve Bayes

# Probability Basics
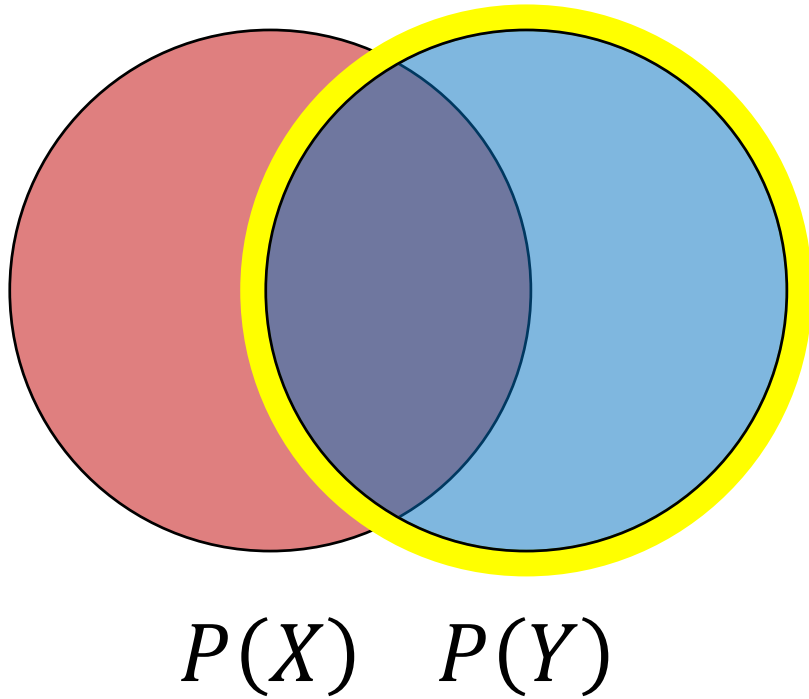


$$P(X)$$

- Single event probability:
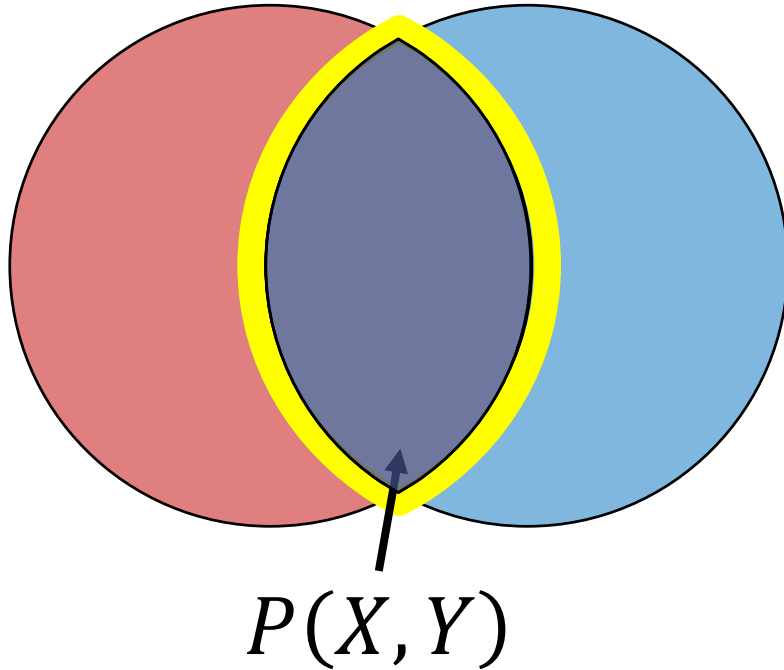
$$P(X)$$

# Probability Basics



$$P(X) \quad P(Y)$$

- Single event probability:

$$P(X), P(Y)$$

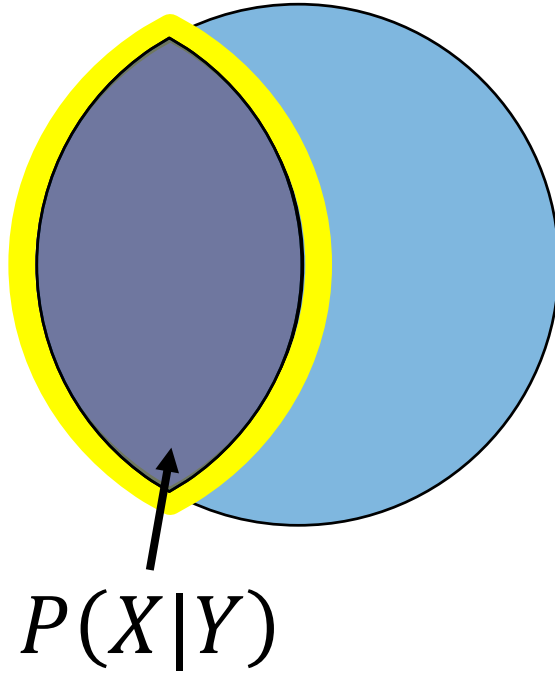# Probability Basics



$$P(X, Y)$$

- Single event probability:

$$P(X), P(Y)$$

- Joint event probability:

$$P(X, Y)$$

# Probability Basics



$P(X|Y)$

- Single event probability:

$$P(X), P(Y)$$

- Joint event probability:

$$P(X, Y)$$

- Conditional probability:

$$P(X|Y)$$

# Probability Basics

$$P(Y|X)$$

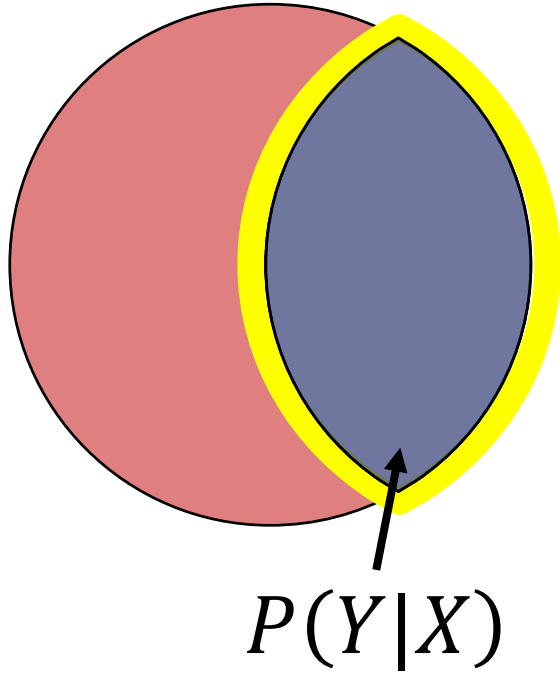- Single event probability:

$$P(X), P(Y)$$

- Joint event probability:

$$P(X, Y)$$
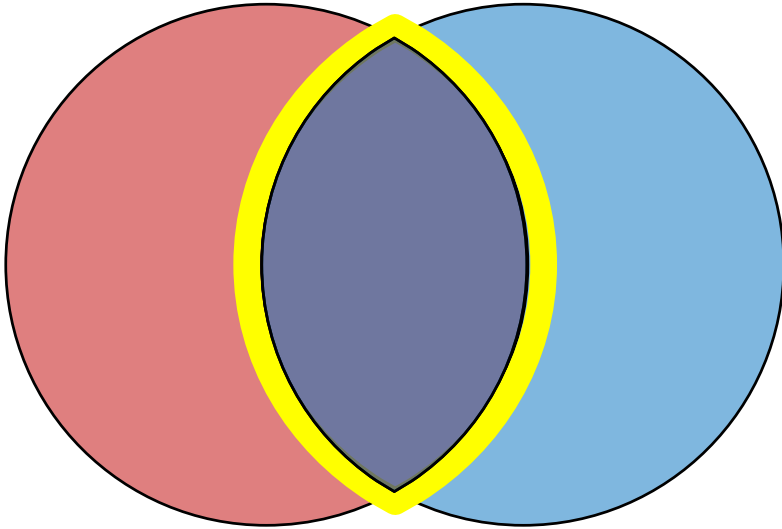
- Conditional probability:

$$P(X|Y), P(Y|X)$$

# Probability Basics



- Single event probability:

$$P(X), P(Y)$$

- Joint event probability:

$$P(X, Y)$$

- Conditional probability:

$$P(X|Y), P(Y|X)$$

- Joint and conditional relationship:

$$P(X, Y) = P(Y|X) * P(X) = P(X|Y) * P(Y)$$

# Bayes Theorem Derivation



- By conditional and joint relationship:

$$P(Y|X) * P(X) = P(X|Y) * P(Y)$$

# Bayes Theorem Derivation



- Use conditional and joint relationship:

$$P(Y|X) * P(X) = P(X|Y) * P(Y)$$

- To invert conditional probability:

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

# Bayes Theorem Derivation



- Use conditional and joint relationship:

$$P(Y|X) * P(X) = P(X|Y) * P(Y)$$

- To invert conditional probability:

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{\boxed{P(X)}}$$

$$P(X) = \sum_Z P(X, Z) = \sum_Z P(X|Z) * P(Z)$$

# Bayes Theorem

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

# Bayes Theorem

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

$$posterior = \frac{likelihood * prior}{evidence}$$

# Naïve Bayes Classification

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

$$posterior = \frac{likelihood * prior}{evidence}$$

(intel) Software

# Training Naïve Bayes

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

**Class** **Feature**

# Training Naïve Bayes: The Naïve Assumption

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

- Difficult to calculate joint probabilities produced by expanding for all features

$$P(C|X) = P(X_1, X_2, \ldots, X_n|C) * P(C)$$
$$P(X_1|X_2, \ldots, X_n, C) * P(X_2, \ldots, X_n|C) * P(C)$$
$$\ldots$$

# Training Naïve Bayes: The Naïve Assumption

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

- **Solution:** assume all features independent of each other

$$P(C|X) = P(X_1|C) * P(X_2|C) * P(X_n|C) * P(C)$$

# Training Naïve Bayes: The Naïve Assumption

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

- **Solution:** assume all features independent of each other

$$P(C|X) = P(X_1|C) * P(X_2|C) * P(X_n|C) * P(C)$$

- This is the "naïve" assumption

$$P(C|X) = P(C) \prod_{i=1}^{n} P(X_i|C)$$

# Training Naïve Bayes

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

- Class assignment is selected based on *maximum a posteriori* (MAP) rule

$$\frac{argmax}{k \in \{1, \ldots K\}} P(C_k) \prod_{i=1}^{n} P(X_i|C_k)$$

# Training Naïve Bayes

- For each class ($C$), calculate probability given features ($X$)

$$P(C|X) = P(X|C) * P(C)$$

- Class assignment is selected based on *maximum a posteriori* (MAP) rule

$$\frac{argmax}{k \in \{1, \ldots K\}} P(C_k) \prod_{i=1}^{n} P(X_i|C_k)$$

Means select potential class with largest value

# The Log Trick

- Multiplying many values together causes computational instability (underflows)

$$\frac{argmax}{k \in \{1, \dots K\}} P(\textcolor{red}{C_k}) \prod_{i=1}^{n} P(\textcolor{blue}{X_i}|\textcolor{red}{C_k})$$

# The Log Trick

- Multiplying many values together causes computational instability (underflows)

$$\frac{argmax}{k \in \{1, \ldots K\}} P(C_k) \prod_{i=1}^{n} P(X_i | C_k)$$

- Work with log values and sum the results

$$\log(P(C_k)) \sum_{i=1}^{n} \log(P(X_i | C_k))$$

# Example: Predicting Tennis With Naïve Bayes

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# Example: Training Naïve Bayes Tennis Model

P(Play=Yes) = 9/14        P(Play=No) = 5/14

Create probability lookup tables based on training data

# Example: Training Naïve Bayes Tennis Model

P(Play=Yes) = 9/14        P(Play=No) = 5/14

| Outlook | Play=Yes | Play=No |
|---------|----------|---------|
| Sunny | 2/9 | 3/5 |
| Overcast | 4/9 | 0/5 |
| Rain | 3/9 | 2/5 |

| Temperature | Play=Yes | Play=No |
|-------------|----------|---------|
| Hot | 2/9 | 2/5 |
| Mild | 4/9 | 2/5 |
| Cool | 3/9 | 1/5 |

Create probability lookup tables based on training data

# Example: Training Naïve Bayes Tennis Model

P(Play=Yes) = 9/14          P(Play=No) = 5/14

| Outlook | Play=Yes | Play=No |
|---------|----------|---------|
| Sunny | 2/9 | 3/5 |
| Overcast | 4/9 | 0/5 |
| Rain | 3/9 | 2/5 |

| Temperature | Play=Yes | Play=No |
|-------------|----------|---------|
| Hot | 2/9 | 2/5 |
| Mild | 4/9 | 2/5 |
| Cool | 3/9 | 1/5 |

| Humidity | Play=Yes | Play=No |
|----------|----------|---------|
| High | 3/9 | 4/5 |
| Normal | 6/9 | 1/5 |

| Wind | Play=Yes | Play=No |
|------|----------|---------|
| Strong | 3/9 | 3/5 |
| Weak | 6/9 | 2/5 |

Create probability lookup tables based on training data

# Example: Predicting Tennis With Naïve Bayes

Predict outcome for the following:

x'=(Outlook=Sunny, Temperature=Cool, Humidity=High, Wind=Strong)

$$P(yes|sunny, cool, high, strong) = P(sunny|yes) * P(cool|yes) *$$
$$P(high|yes)*P(strong|yes)*P(yes)$$

$$P(no|sunny, cool, high, strong) = P(sunny|no) * P(cool|no) *$$
$$P(high|no)*P(strong|no)*P(no)$$

# Example: Predicting Tennis With Naïve Bayes

Predict outcome for the following:

x'=(Outlook=Sunny, Temperature=Cool, Humidity=High, Wind=Strong)

| Feature | Play=Yes | Play=No |
|---|---|---|
| Outlook=Sunny | 2/9 | 3/5 |

# Example: Predicting Tennis With Naïve Bayes

Predict outcome for the following:

x'=(Outlook=Sunny, Temperature=Cool, Humidity=High, Wind=Strong)

| Feature | Play=Yes | Play=No |
|---|---|---|
| Outlook=Sunny | 2/9 | 3/5 |
| Temperature=Cool | 3/9 | 1/5 |
| Humidity=High | 3/9 | 4/5 |
| Wind=Strong | 3/9 | 3/5 |
| **Overall Label** | **9/14** | **5/14** |

# Example: Predicting Tennis With Naïve Bayes

Predict outcome for the following:

x'=(Outlook=Sunny, Temperature=Cool, Humidity=High, Wind=Strong)

| Feature | Play=Yes | Play=No |
|---|---|---|
| Outlook=Sunny | 2/9 | 3/5 |
| Temperature=Cool | 3/9 | 1/5 |
| Humidity=High | 3/9 | 4/5 |
| Wind=Strong | 3/9 | 3/5 |
| **Overall Label** | **9/14** | **5/14** |
| **Probability** | **0.0053** | **0.0206** |

# Example: Predicting Tennis With Naïve Bayes

Predict outcome for the following:

x'=(Outlook=Sunny, Temperature=Cool, Humidity=High, Wind=Strong)

| Feature | Play=Yes | Play=No |
|---|---|---|
| Outlook=Sunny | 2/9 | 3/5 |
| Temperature=Cool | 3/9 | 1/5 |
| Humidity=High | 3/9 | 4/5 |
| Wind=Strong | 3/9 | 3/5 |
| **Overall Label** | **9/14** | **5/14** |
| **Probability** | **0.0053** | **0.0206** |

# Laplace Smoothing

- **Problem:** categories with no entries result in a value of "0" for conditional probability

$$P(C|X) = P(X_1|C) * P(X_2|C) * P(C)$$

# Laplace Smoothing

- **Problem:** categories with no entries result in a value of "0" for conditional probability

$$P(C|X) = \boxed{P(X_1|C)} * P(X_2|C) * P(C)$$

0

# Laplace Smoothing

0

- **Problem:** categories with no entries result in a value of "0" for conditional probability

$$P(C|X) = \boxed{P(X_1|C)} * P(X_2|C) * P(C)$$

- **Solution:** add "1" to numerator and denominator of empty categories

$$P(X_1|C) = \frac{1}{Count(C) + n}$$

$$P(X_2|C) = \frac{Count(X_2 \ \& \ C) + 1}{Count(C) + m}$$

# Types of Naïve Bayes

| Naïve Bayes Model | Data Type |
|---|---|
| Bernoulli | Binary (T/F) |

# Types of Naïve Bayes

| Naïve Bayes Model | Data Type |
|---|---|
| Bernoulli | Binary (T/F) |
| Multinomial | Discrete (e.g. count) |

# Types of Naïve Bayes

| Naïve Bayes Model | Data Type |
|---|---|
| Bernoulli | Binary (T/F) |
| Multinomial | Discrete (e.g. count) |
| Gaussian | Continuous |

# Combining Feature Types

| Problem | • Model features contain different data types (continuous and categorical) |

# Combining Feature Types

**Problem**

- Model features contain different data types (continuous and categorical)

**Solutions**

- **Option 1:** Bin continuous features to create categorical ones and fit multinomial model

# Combining Feature Types

**Problem**

- Model features contain different data types (continuous and categorical)

**Solutions**

- **Option 1:** Bin continuous features to create categorical ones and fit multinomial model

- **Option 2:** Fit Gaussian model on continuous features and multinomial on categorical features; combine to create "meta model" (week 10)

# Distributed Computing with Naïve Bayes

- Well-suited for large data and distributed computing—limited parameters and log probabilities are a summation

- Scikit-Learn implementations contain a "partial_fit" method designed for out-of-core calculations

# Naïve Bayes: The Syntax

**Import the class containing the classification method**

from sklearn.naive_bayes import **BernoulliNB**

# Naïve Bayes: The Syntax

**Import the class containing the classification method**

from sklearn.naive_bayes import **BernoulliNB**

**Create an instance of the class**
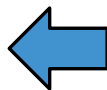
**BNB** = **BernoulliNB**(alpha=1.0)

# Naïve Bayes: The Syntax

**Import the class containing the classification method**

from sklearn.naive_bayes import **BernoulliNB**

**Create an instance of the class**

**BNB** = **BernoulliNB**(alpha=1.0)

Laplace
smoothing
parameter

# Naïve Bayes: The Syntax

**Import the class containing the classification method**

 from sklearn.naive_bayes import **BernoulliNB**

**Create an instance of the class**

 **BNB** = **BernoulliNB**(alpha=1.0)

**Fit the instance on the data and then predict the expected value**

 **BNB** = **BNB**.**fit**(X_train, y_train)

 y_predict = **BNB**.**predict**(X_test)

# Naïve Bayes: The Syntax

**Import the class containing the classification method**

    **from sklearn.naive_bayes import BernoulliNB**

**Create an instance of the class**

    **BNB = BernoulliNB(alpha=1.0)**

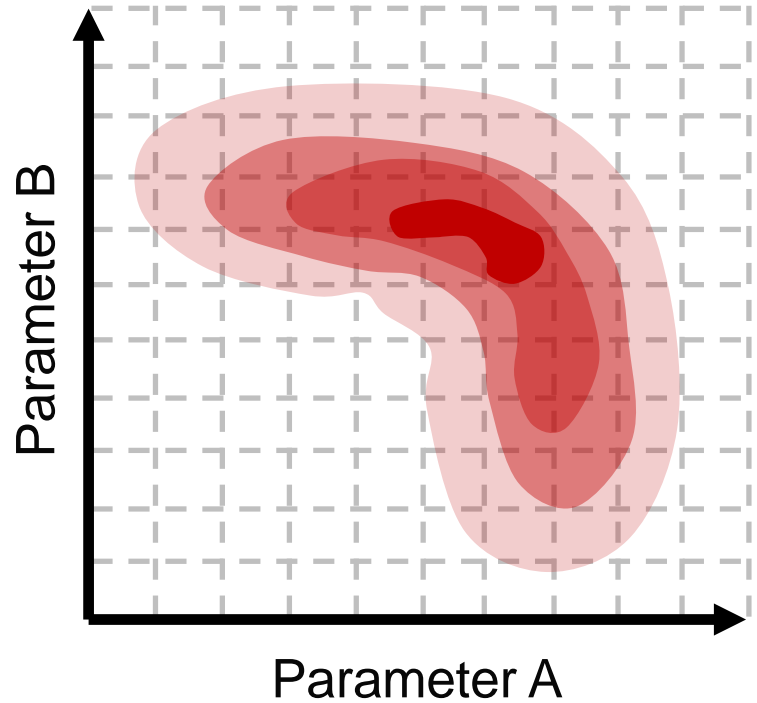**Fit the instance on the data and then predict the expected value**

    **BNB = BNB.fit(X_train, y_train)**

    **y_predict = BNB.predict(X_test)**

**Other naïve Bayes models: MultinomialNB, GaussianNB.**

# Grid Search & Pipelines
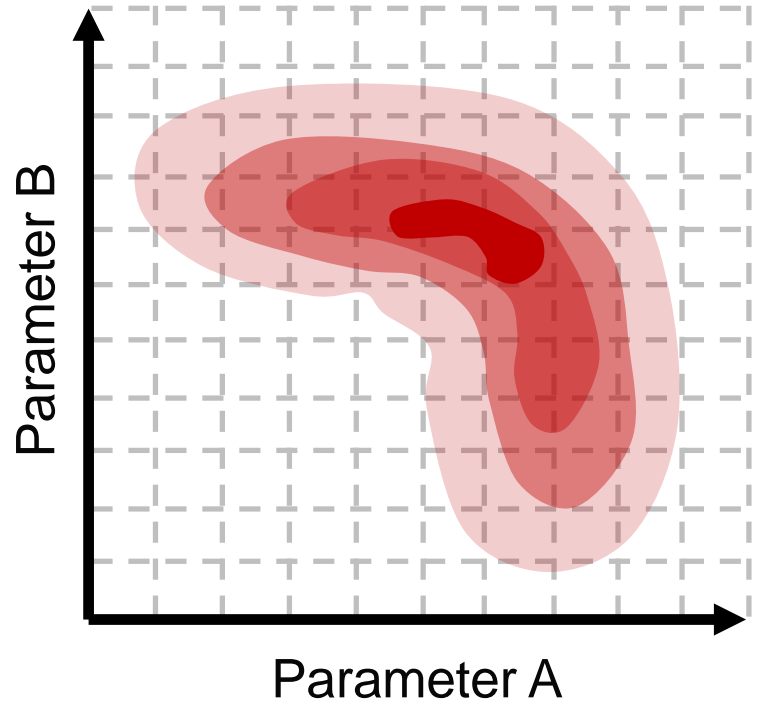
# Generalized Hyperparameter Grid Search

- Hyperparameter selection for regularization / better models requires cross validation on training data

- Linear and logistic regression methods have classes devoted to grid search (e.g. LassoCV)



Parameter B

Parameter A

# Generalized Hyperparameter Grid Search

- Grid search can be useful for other methods too, so a generalized method is desirable

- Scikit-learn contains GridSearchCV, which performs a grid search with parameters using cross validation

Parameter B

Parameter A

# Grid Search with Cross Validation: The Syntax

**Import the class containing the grid search method**

from sklearn.linear_model import **LogisticRegression**

from sklearn.model_selection import **GridSearchCV**

# Grid Search with Cross Validation: The Syntax

**Import the class containing the grid search method**

from sklearn.linear_model import **LogisticRegression**

from sklearn.model_selection import **GridSearchCV**

**Create an instance of the estimator and grid search class**

**LR** = **LogisticRegression**(penalty='l2')

**GS** = **GridSearchCV**(**LR**, param_grid={'c':[0.001, 0.01, 0.1]},

scoring='accuracy', cv=4)

# Grid Search with Cross Validation: The Syntax

**Import the class containing the grid search method**

from sklearn.linear_model import **LogisticRegression**

from sklearn.model_selection import **GridSearchCV**

**Create an instance of the estimator and grid search class**

logistic regression method

**LR** = **LogisticRegression**(penalty='l2')

**GS** = **GridSearchCV**(**LR**, param_grid={'c':[0.001, 0.01, 0.1]},

scoring='accuracy', cv=4)

# Grid Search with Cross Validation: The Syntax

**Import the class containing the grid search method**

from sklearn.linear_model import **LogisticRegression**

from sklearn.model_selection import **GridSearchCV**

**Create an instance of the estimator and grid search class**

**LR** = **LogisticRegression**(penalty='l2')

**GS** = **GridSearchCV**(**LR**, param_grid={'c':[0.001, 0.01, 0.1]},

scoring='accuracy', cv=4)

**Fit the instance on the data to find the best model and then predict**

**GS** = **GS**.**fit**(X_train, y_train)

y_train = **GS**.**predict**(X_test)

# Optimizing the Rest of the Pipeline

- Grid searches enable model parameters to be optimized

# Optimizing the Rest of the Pipeline

- Grid searches enable model parameters to be optimized

- How can this be incorporated with other steps of the process (e.g. feature extraction and transformation)?

# Optimizing the Rest of the Pipeline

- Grid searches enable model parameters to be optimized

- How can this be incorporated with other steps of the process (e.g. feature extraction and transformation)?
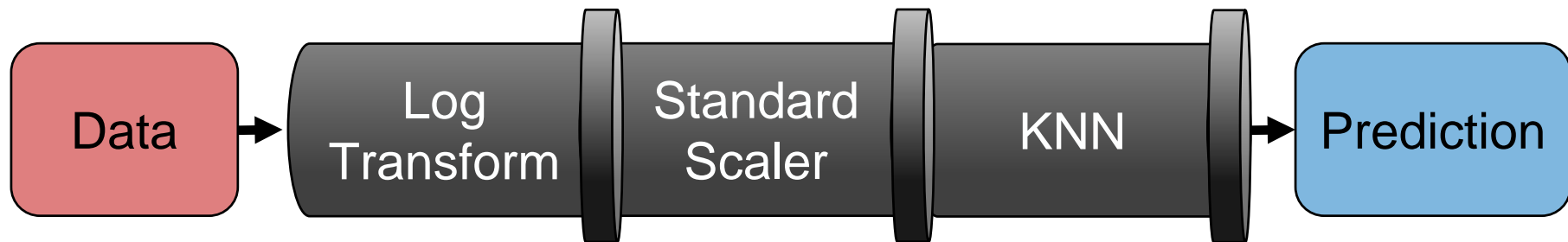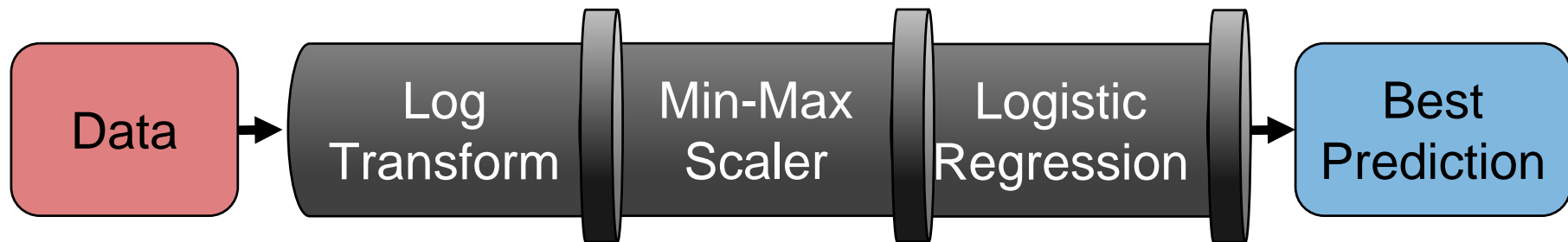


Pipelines!

# Automating Machine Learning with Pipelines

- Machine learning models often selected empirically

# Automating Machine Learning with Pipelines

- Machine learning models often selected empirically
- By trying different processing methods and tuning multiple models
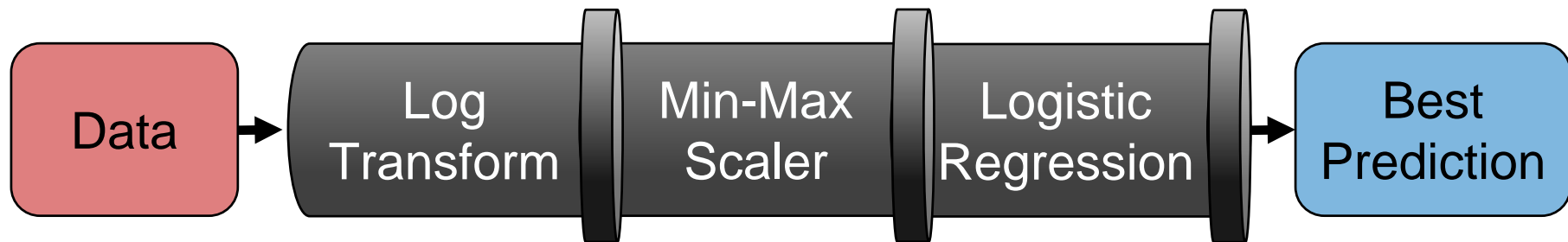
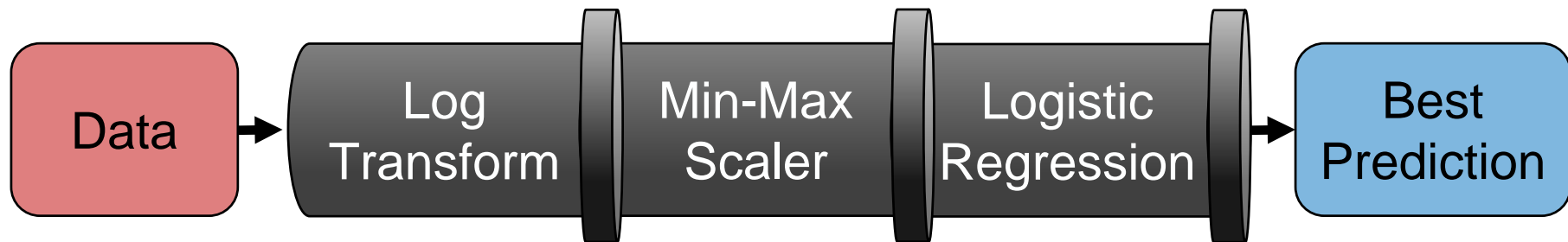# Automating Machine Learning with Pipelines

- Machine learning models often selected empirically
- By trying different processing methods and tuning multiple models
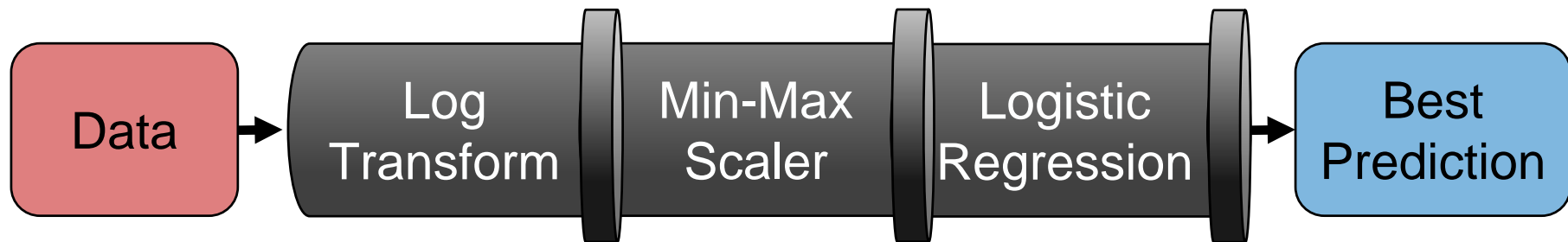


How to automate this process?

# Automating Machine Learning with Pipelines

- Pipelines in Scikit-Learn allow feature transformation steps and models to be chained together

# Automating Machine Learning with Pipelines

- Pipelines in Scikit-Learn allow feature transformation steps and models to be chained together
- Successive steps perform 'fit' and 'transform' before sending data to the next step
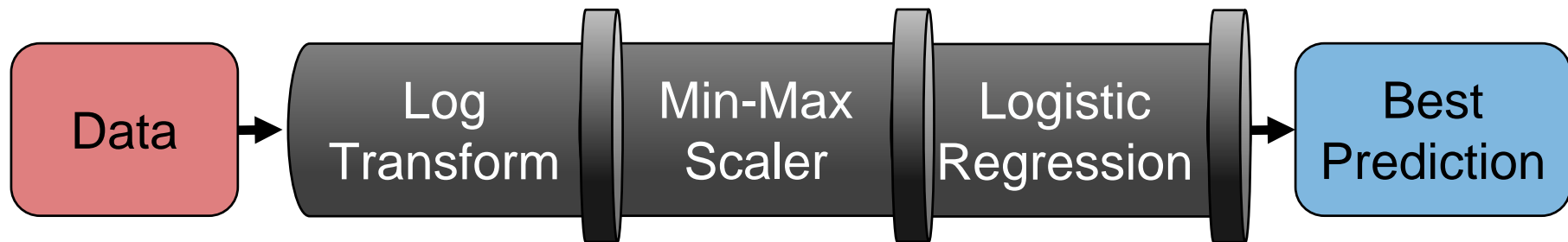
# Automating Machine Learning with Pipelines

- Pipelines in Scikit-Learn allow feature transformation steps and models to be chained together
- Successive steps perform 'fit' and 'transform' before sending data to the next step

Data → | Log Transform | Min-Max Scaler | Logistic Regression | → Best Prediction

Pipelines make automation and reproducibility easier!

# Pipelines: The Syntax

**Import the class containing the pipeline method**

    **from sklearn.pipeline import Pipeline**

# Pipelines: The Syntax

**Import the class containing the pipeline method**

```
from sklearn.pipeline import Pipeline
```

**Create an instance of the class with estimators**

```
estimators = [('scaler', MinMaxScaler()), ('lasso', Lasso())]
Pipe = Pipeline(estimators)
```

# Pipelines: The Syntax

**Import the class containing the pipeline method**

> **from sklearn.pipeline import Pipeline**

**Create an instance of the class with estimators**

feature scaler class

**estimators = [('scaler', MinMaxScaler()), ('lasso', Lasso())]**

> **Pipe = Pipeline(estimators)**

# Pipelines: The Syntax

**Import the class containing the pipeline method**

**from sklearn.pipeline import Pipeline**

**Create an instance of the class with estimators**

**estimators = [('scaler', MinMaxScaler()), ('lasso', Lasso())]**

**Pipe = Pipeline(estimators)**

lasso model class

# Pipelines: The Syntax

**Import the class containing the pipeline method**

    from sklearn.pipeline import **Pipeline**

**Create an instance of the class with estimators**

  estimators = [('scaler', **MinMaxScaler()**), ('lasso', **Lasso()**)]

    **Pipe** = **Pipeline**(estimators)

**Fit the instance on the data and then predict the expected value**

    **Pipe** = **Pipe**.**fit**(X_train, y_train)

    y_predict = **Pipe**.**predict**(X_test)

# Pipelines: The Syntax

**Import the class containing the pipeline method**

    **from sklearn.pipeline import Pipeline**

**Create an instance of the class with estimators**

  **estimators = [('scaler', MinMaxScaler()), ('lasso', Lasso())]**

    **Pipe = Pipeline(estimators)**

**Fit the instance on the data and then predict the expected value**

    **Pipe = Pipe.fit(X_train, y_train)**

    **y_predict = Pipe.predict(X_test)**

**Features can be combined from different transform method using FeatureUnion**

# Legal Notices and Disclaimers

This presentation is for informational purposes only. INTEL MAKES NO WARRANTIES, EXPRESS OR IMPLIED, IN THIS SUMMARY.

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at intel.com.

This sample source code is released under the Intel Sample Source Code License Agreement.

Intel and the Intel logo are trademarks of Intel Corporation in the U.S. and/or other countries.

*Other names and brands may be claimed as the property of others.