



## #3 Intro to ML

---

AI Training

# The Machine Learning Landscape

# What is Machine Learning?

Machine Learning is the science (and art) of programming computers so they can learn from data.

Here is a slightly more general definition:

[Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed.

—Arthur Samuel, 1959

# Example: Spam Filter as a ML Program

- **Goal:** Automatically detect spam emails
- **Input:**
  - Labeled spam emails (e.g., flagged by users)
  - Labeled non-spam (ham) emails
- **Output:**
  - A model that learns to flag future spam emails

# Why Use Machine Learning?

Machine Learning is great for:

- **Rule-heavy problems**

When traditional solutions need constant fine-tuning or many rules, ML can simplify and outperform them.

```
if flat.price is <= 10.0:  
    then send offer;  
  
if flat is big...  
if flat is new  
...
```

```
try to find good flats:  
adapt model to learn from err  
repeat until threshold is ;
```

- **Complex or unsolved problems**

Where traditional methods fail, ML can uncover patterns and find workable solutions.

- **Changing environments**

ML models adapt as new data arrives – ideal for dynamic or evolving systems.

# Examples of Applications

- Analyzing images of products on a production line to automatically classify them
- Detecting tumors in brain scans
- Automatically classifying news articles
- Automatically flagging offensive comments on discussion forums
- Summarizing long documents automatically
- Creating a chatbot or a personal assistant
- Forecasting your company's revenue next year, based on many performance metrics
- Making your app react to voice commands
- Detecting credit card fraud
- Segmenting clients based on their purchases so that you can design a different marketing strategy for each segment
- Representing a complex, high-dimensional dataset in a clear and insightful diagram
- Recommending a product that a client may be interested in, based on past purchases
- Building an intelligent bot for a game

# Types of Machine Learning Systems

# Three types of machine learning

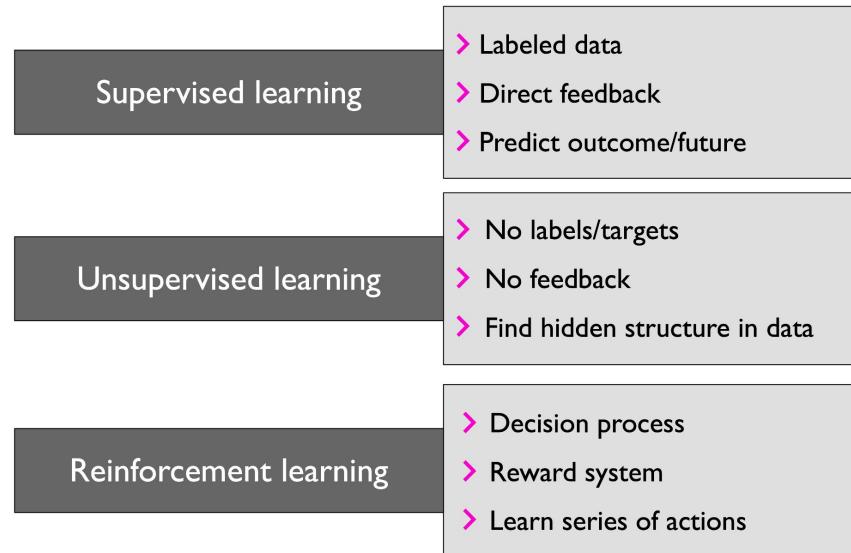


Figure 1.1: The three different types of machine learning

# Supervised Learning

- **Goal:** Learn a model from **labeled training data** to make predictions on **unseen/future data**
- **Why “Supervised”?**
  - Each training example comes with a **known label** (the “correct answer”)
- **What it does:**
  - Models the relationship between **inputs** (features) and **outputs** (labels)
- **Another way to think about it:**
  - Supervised learning = “**learning from labels**”

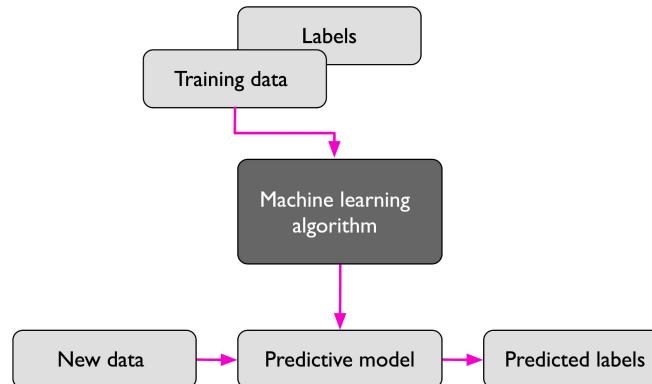


Figure 1.2: Supervised learning process

# Supervised Learning: Classification

- **Goal:** Predict categorical class labels for new data points
- **Labels are discrete, unordered values** representing group membership

## Example: Binary Classification

- Task: **Spam detection**
- Two possible classes: **Spam vs. Not Spam**
- Model learns rules from labeled examples

## Example: Multiclass Classification

- Task: **Handwritten character recognition**
- Classes: Letters A-Z
- Model predicts the correct letter from input
- Limitation: Can't recognize digits (0-9) unless seen during training

# Supervised Learning: Classification

- Scenario: 30 training examples
  - 15 labeled Class A, 15 labeled Class B
- Goal:
  - Use a supervised ML algorithm to learn a **decision boundary**
  - This boundary (dashed line) separates Class A from Class B

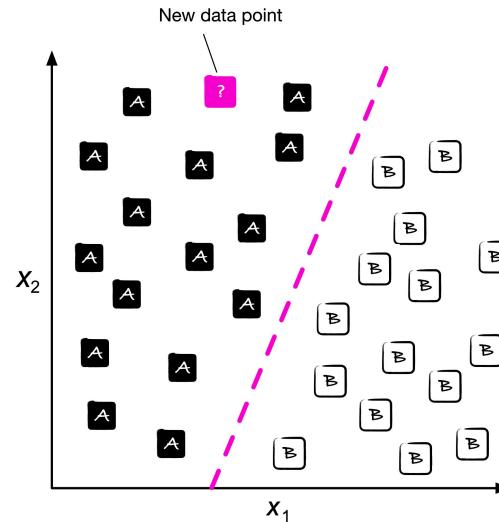


Figure 1.3: Classifying a new data point

# Supervised Learning: Regression

- Unlike classification, regression predicts continuous values, not categories
- Also known as **regression analysis**

Example: Predicting a student's SAT Math Scores

- Feature: Time spent studying
- Target: Final score

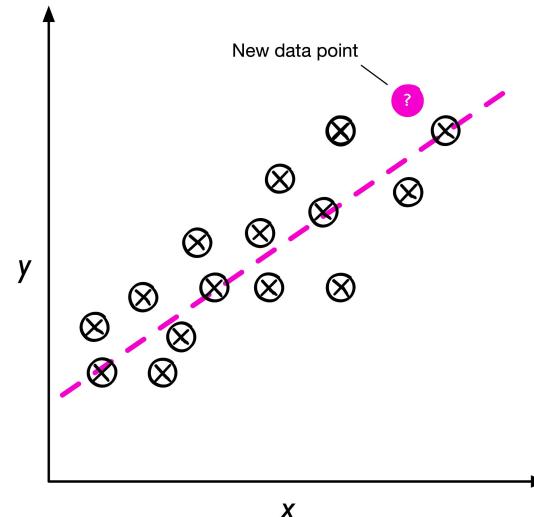


Figure 1.4: A linear regression example

# Unsupervised Learning

- Works with **unlabeled data**
- No predefined **outcomes or target variables**
- Goal: **Discover hidden patterns or underlying structure** in the data
- No external supervision or feedback (like rewards or labels)

## Examples of Unsupervised Learning Tasks:

- **Clustering:** Grouping similar data points
- **Dimensionality reduction:** Simplifying data while preserving structure

# Unsupervised Learning: Clustering

- Clustering is an exploratory technique used to **organize data into meaningful groups (clusters)**
- No prior knowledge of group memberships required
- Groups are formed based on **similarity within clusters and dissimilarity between clusters**

Why It's Useful:

- Helps **reveal hidden patterns** or relationships
- Also known as **unsupervised classification**

Example:

- **Marketing:** Identify distinct customer segments based on interests to tailor marketing strategies

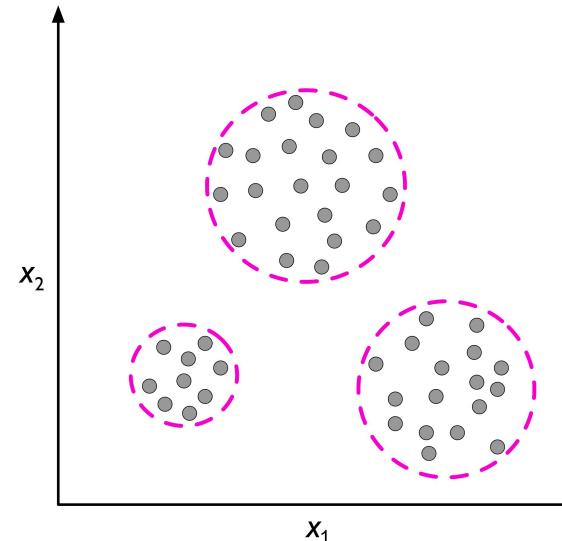


Figure 1.5: How clustering works

# Reinforcement Learning

- An agent learns by interacting with an **environment**
- Learns to take actions that **maximize cumulative reward** over time
- Feedback comes as a **reward signal**, not a correct label

## How it Works

- Agent observes the current **state**
- Takes an **action**
- Receives a **reward** based on that action
- Learns through **trial-and-error or planning**

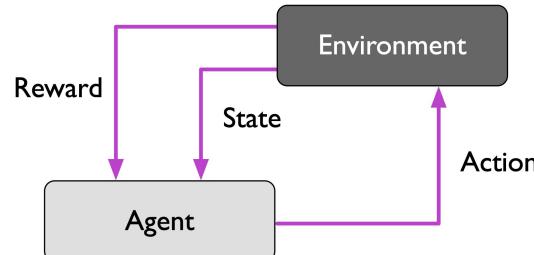


Figure 1.6: Reinforcement learning process

# Reinforcement Learning

## Key Differences from Supervised Learning

- No ground truth label for each input
- Feedback is **delayed and indirect** (via reward function)
- Goal: Learn a **policy** that maps states to optimal actions

## Example: Chess AI

- **State:** Current board position
- **Action:** Chosen move
- **Reward:** +1 for win, 0 for draw, -1 for loss
- Learns strategies by playing many games

# Machine Learning Terminology

## Training Example

- A row in a dataset (a single data point)
- Also known as: **observation, record, instance, or sample**

## Training

- The process of **model fitting**
- In parametric models: similar to **parameter estimation**

## Feature (x)

- A **column in the dataset** (an input variable)
- Also known as: **predictor, input, variable, attribute, or covariate**

## Target (y)

- The value to be predicted (the output variable)
- Also known as: **outcome, label, response, dependent variable, ground truth**

## Loss Function / Cost Function

- Measures the **error** between predicted and actual values
- Synonyms: **error function, objective function**

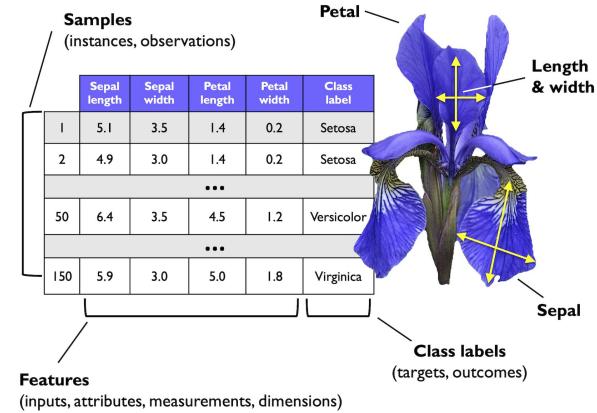


Figure 1.7: The Iris dataset

# Main Challenges of Machine Learning

# Machine Learning Challenges

Since your task is to select a learning algorithm and train it on some data, **two key things can go wrong**:

## Bad Data

- Insufficient data quantity
- Data that does not represent the real-world problem
- Poor-Quality Data
- Irrelevant or redundant features

## Bad Algorithm

- Wrong choice of model for the task
- Overfitting or underfitting
- Poor hyperparameter tuning

*Even the best algorithm can fail on poor data – and great data can be wasted on the wrong model.*

*Bad Data*

# Insufficient Quantity of Training Data

## Human Learning vs. Machine Learning

For a toddler to learn what an apple is, all it takes is for you to point to an apple and say “apple” (possibly repeating this procedure a few times). Now the child is able to recognize apples in all sorts of colors and shapes. Genius. Machine Learning is not quite there yet;

## ML Needs Lots of Data

- **Thousands** of examples for simple tasks
- **Millions** for complex tasks (e.g., image or speech recognition - unless you can reuse parts of an existing model known as *transfer learning*)
- More data = better generalization, less overfitting

# Nonrepresentative Training Data

In order to generalize well, it is crucial that your training data be representative of the new cases you want to generalize to.

## Common Pitfalls

- **Sampling bias:** Data collected from one group or condition but applied to another
- **Temporal drift:** Training data is outdated and no longer reflects the current environment
- **Coverage gaps:** Important cases or variations are missing from training examples
- **Sampling noise:** Random variations in the sample that cause it to misrepresent the true distribution of the population (e.g., over- or under-representing certain groups or outcomes)

# Poor-Quality Data

- A major part of a data scientist's time is spent **cleaning and preprocessing data**
- High-quality data often matters **more than complex algorithms**

## Garbage In, Garbage Out

- Errors, noise, and outliers in training data can **hide real patterns**
- Leads to poor generalization and weak model performance

## Examples of Data Issues:

- **Outliers:** Remove or manually correct obvious errors
- **Missing values:**
  - Drop the feature or the affected instances
  - **Impute** missing values (e.g., use median or average)
  - Train **separate models** with and without the feature

# Irrelevant Features

## Garbage In, Garbage Out

- Your model's success depends on having **relevant, meaningful features**
- Too many irrelevant features = noise, overfitting, and poor performance

## The Art of Feature Engineering

A key step in building effective ML models. It includes:

- **Feature selection:** Choosing the most useful features from existing ones
- **Feature extraction:** Creating better features by combining or transforming existing ones
  - (e.g., via dimensionality reduction techniques like PCA)
- **Feature creation:** Gathering new data or constructing new features from domain knowledge

*Bad Algorithm*

# Overfitting the Training Data

- Happens when a model learns the **training data** too well, including its **noise and outliers**
- Performs well on training data but **poorly on new, unseen data**

## Example:

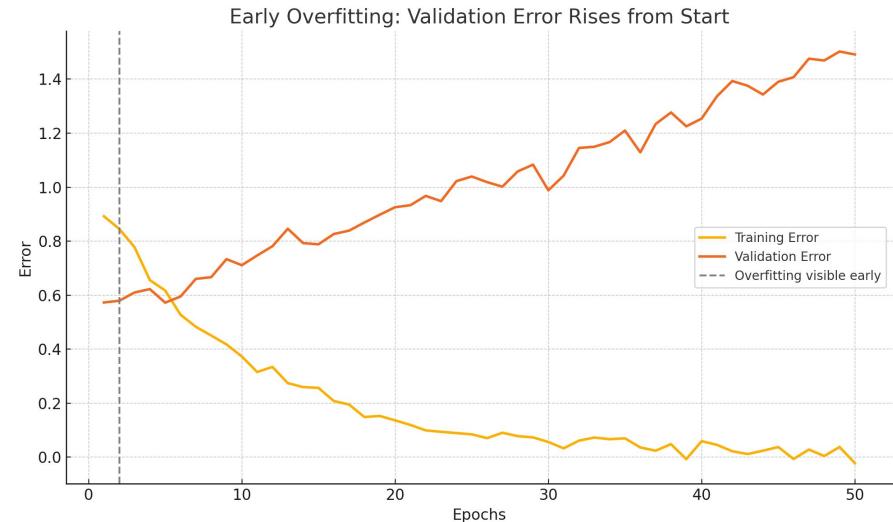
- A complex model fits every fluctuation in training data, but **fails to generalize**  
*(like memorizing answers instead of understanding the material)*

## How to Fix Overfitting

- Use a **simpler model** (fewer parameters)
- **Regularize** the model (e.g., L1 or L2 regularization) - to reduce the model complexity
- Use more **training data** if possible
- Apply **cross-validation** to monitor generalization
- Use techniques like **early stopping** or **dropout** (in neural networks)

# Signs of Overfitting! Validation set

- Detecting overfitting is only possible once we move to the testing phase.
- Don't trust yourself
- Overfitting is massive for smaller datasets
- Ideally have a part of the dataset you don't have access to (validation set)
- Even some signs for \*huge\* datasets (imagenet)



# Underfitting the Training Data

- Occurs when a model is **too simple** to capture the underlying patterns in the data
- Leads to **poor performance**, even on the training set

Example:

- A **linear model** trying to predict life satisfaction – too simplistic for a complex reality

How to Fix Underfitting

- Select a **more powerful model**, with more parameters.
- Feed better features to the learning algorithm (**feature engineering**).
- **Reduce the constraints** on the model (e.g., reduce the regularization hyperparameter).

# Underfitting vs Overfitting

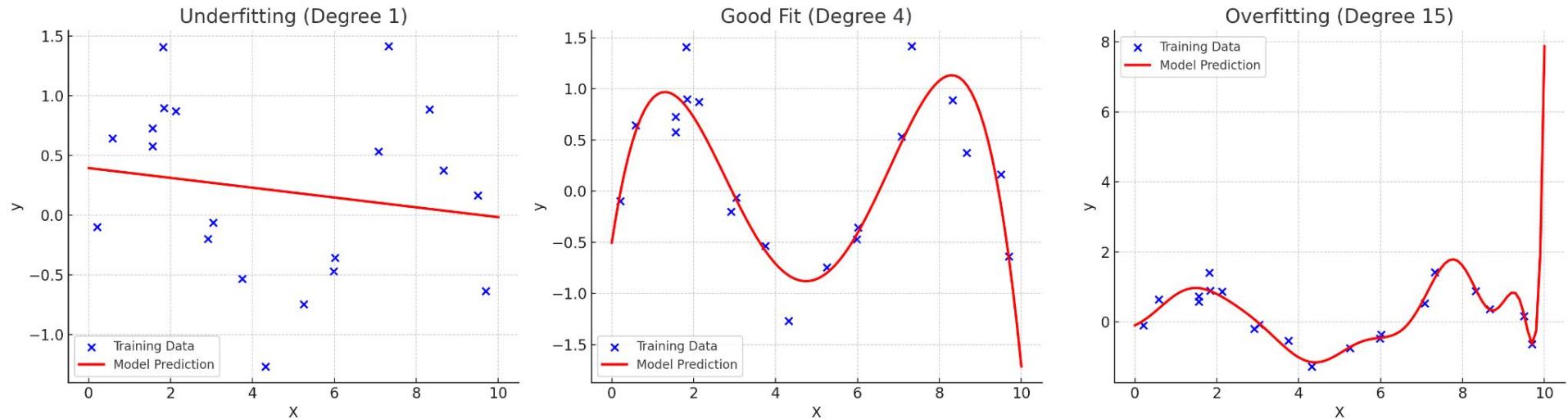


Figure 1.9: Underfitting vs Overfitting

# Testing and Validating

# Testing and Validating your Model

The only way to know how well a model will generalize to new cases is to actually try it out on new cases.

## Goal: Estimate Generalization Performance

- We want to know how well the model performs on **unseen data**
- The safest way: **test it on data it hasn't seen during training**

## Train/Test Split

- **Training Set:** Used to train the model
- **Test Set:** Used to evaluate how well the model generalizes
- The error on the test set is called the **generalization error** (or **out-of-sample error**)
- **Common Practice:** 80% training, 20% testing
  - For very large datasets, even a 1% test set can be enough for reliable evaluation

## Overfitting Warning

- If **training error is low but test error is high** → your model is overfitting

# Testing and Validating your Model

The Problem:

- Choosing the best model and hyperparameters by testing repeatedly on the test set can lead to **overfitting to the test set**
- This makes your final model **less reliable** on truly unseen data

The Solution: Holdout Validation

1. **Split training data into:**
  - A training subset
  - A validation set (aka dev set)
2. **Use the validation set to:**
  - Compare different models
  - Tune hyperparameters (e.g., regularization strength)
3. **Evaluate only once on the test set** for a true estimate of generalization error

# A roadmap for building machine learning systems

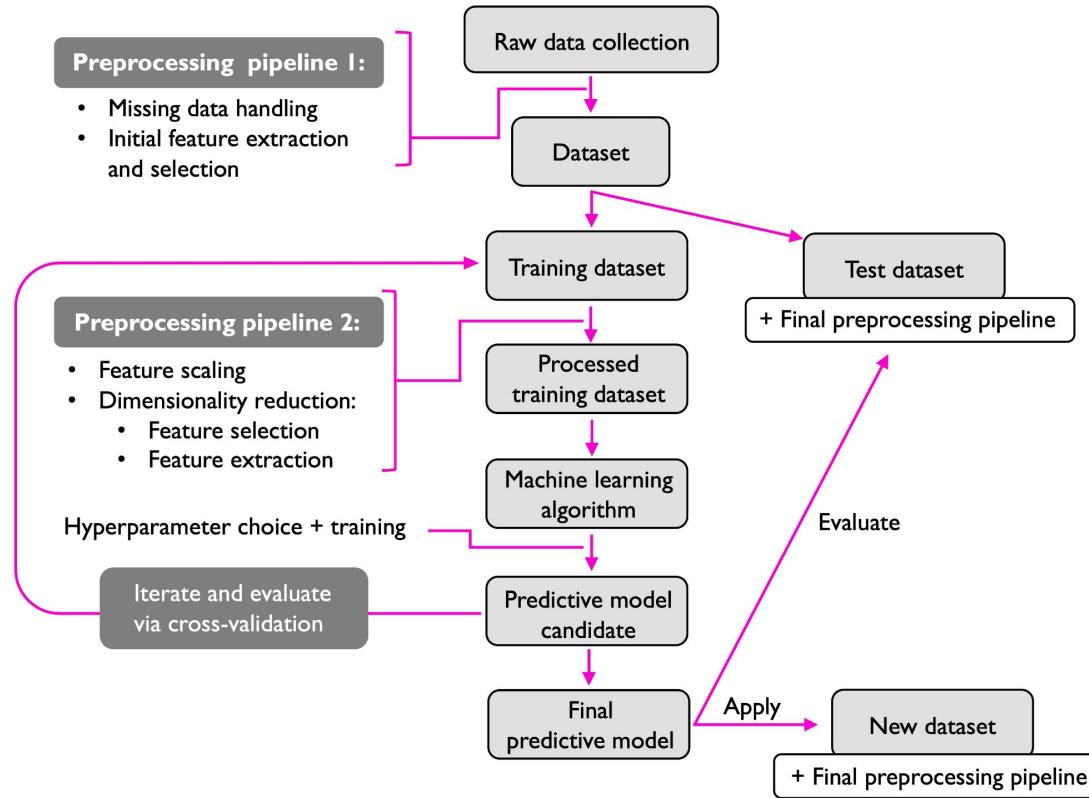


Figure 1.10: Predictive modeling workflow

# Packages for Machine Learning

- NumPy
- SciPy
- Scikit-learn
- Matplotlib
- Pandas

*Check Your Understanding*

# Training and Test Sets

Explore the options below.

We looked at a process of using a test set and a training set to drive iterations of model development. On each iteration, we'd train on the training data and evaluate on the test data, using the evaluation results on test data to guide choices of and changes to various model hyperparameters like learning rate and features. Is there anything wrong with this approach? (Pick only one answer.)

- A. This is computationally inefficient. We should just pick a default set of hyperparameters and live with them to save resources.
- B. Totally fine, we're training on training data and evaluating on separate, held-out test data.
- C. Doing many rounds of this procedure might cause us to implicitly fit to the peculiarities of our specific test set.

# Training and Test Sets

Explore the options below.

We looked at a process of using a test set and a training set to drive iterations of model development. On each iteration, we'd train on the training data and evaluate on the test data, using the evaluation results on test data to guide choices of and changes to various model hyperparameters like learning rate and features. Is there anything wrong with this approach? (Pick only one answer.)

- A. This is computationally inefficient. We should just pick a default set of hyperparameters and live with them to save resources.
- B. Totally fine, we're training on training data and evaluating on separate, held-out test data.
- C. Doing many rounds of this procedure might cause us to implicitly fit to the peculiarities of our specific test set.

# Supervised Learning with Scikit-Learn

# Recap: Supervised Learning

- The predicted values are known
- Aim: Predict the target values of unseen data, given the features

	Features					Target variable
	points_per_game	assists_per_game	rebounds_per_game	steals_per_game	blocks_per_game	position
0	26.9	6.6	4.5	1.1	0.4	Point Guard
1	13	1.7	4	0.4	1.3	Center
2	17.6	2.3	7.9	1.00	0.8	Power Forward
3	22.6	4.5	4.4	1.2	0.4	Shooting Guard

# Before you use Supervised Learning

- Requirements:
  - No missing values
  - Data in numeric format
  - Data stored in pandas DataFrame or NumPy array
- Perform Exploratory Data Analysis (EDA) first

# scikit-learn syntax

```
from sklearn.module import Model  
model = Model()  
model.fit(X, y)  
predictions = model.predict(X_new)  
print(predictions)
```

```
array([0, 0, 0, 0, 1, 0])
```

# Classification

# Classifying labels of unseen data

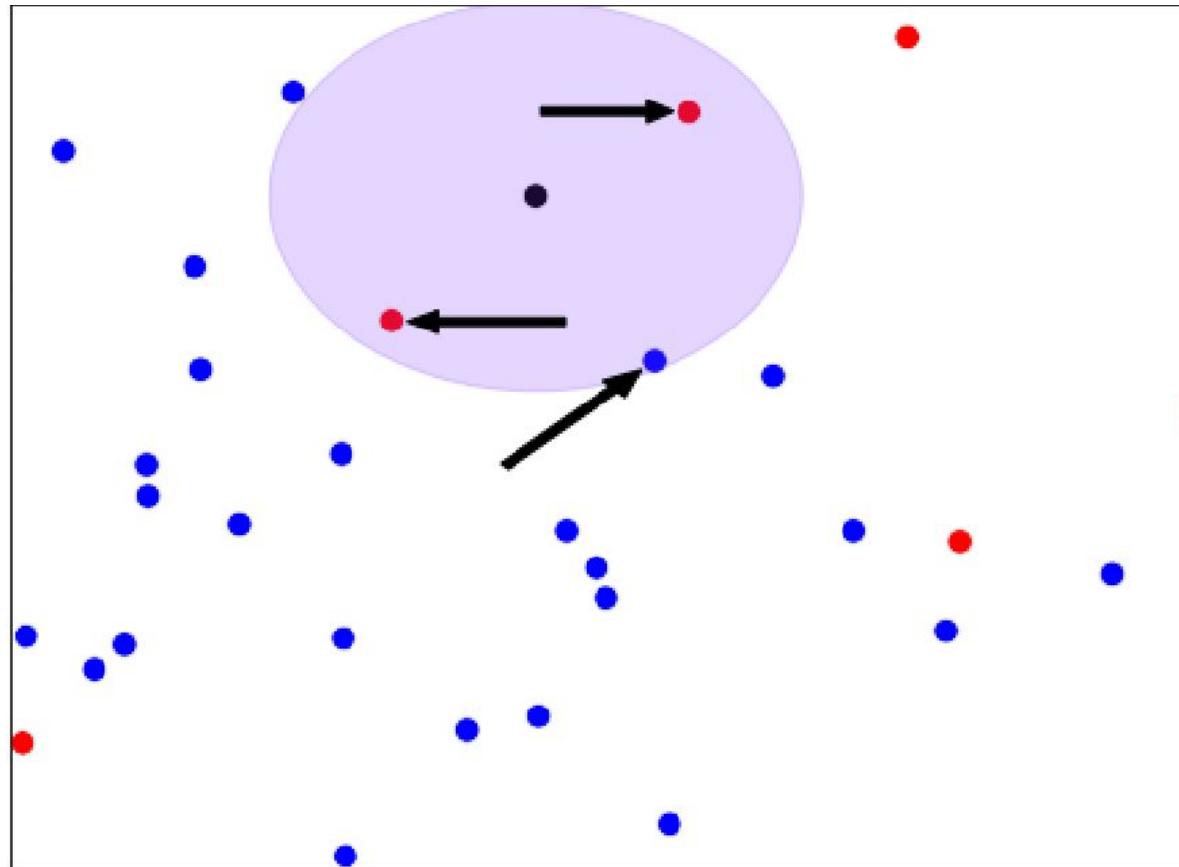
1. Build a model
  2. Model learns from the labeled data we pass to it
  3. Pass unlabeled data to the model as input
  4. Model predicts the labels of the unseen data
- 
- Labeled data = training data

# k-Nearest Neighbors Algorithm

- Predict the label of a data point by
  - Looking at the  $k$  closest labeled data points
  - Taking a majority vote

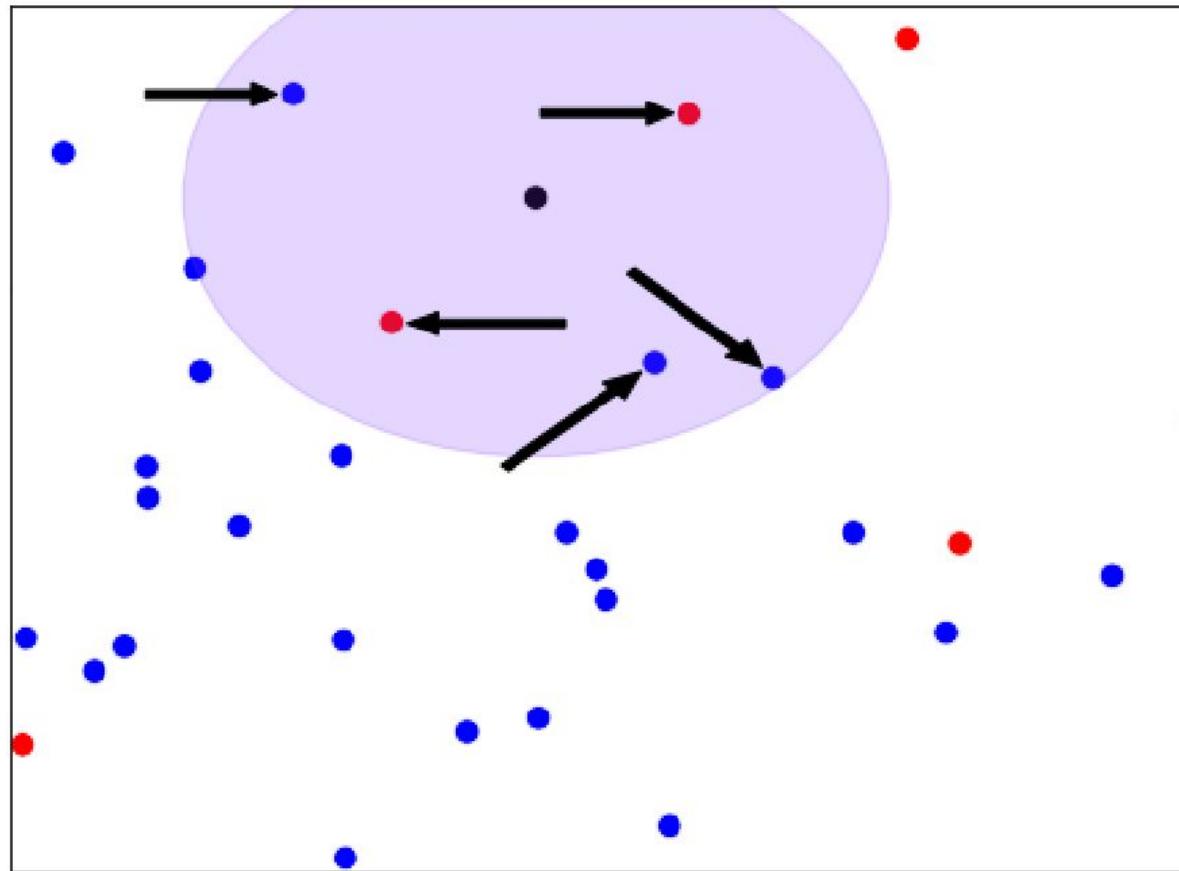
# k-Nearest Neighbors

- K=3

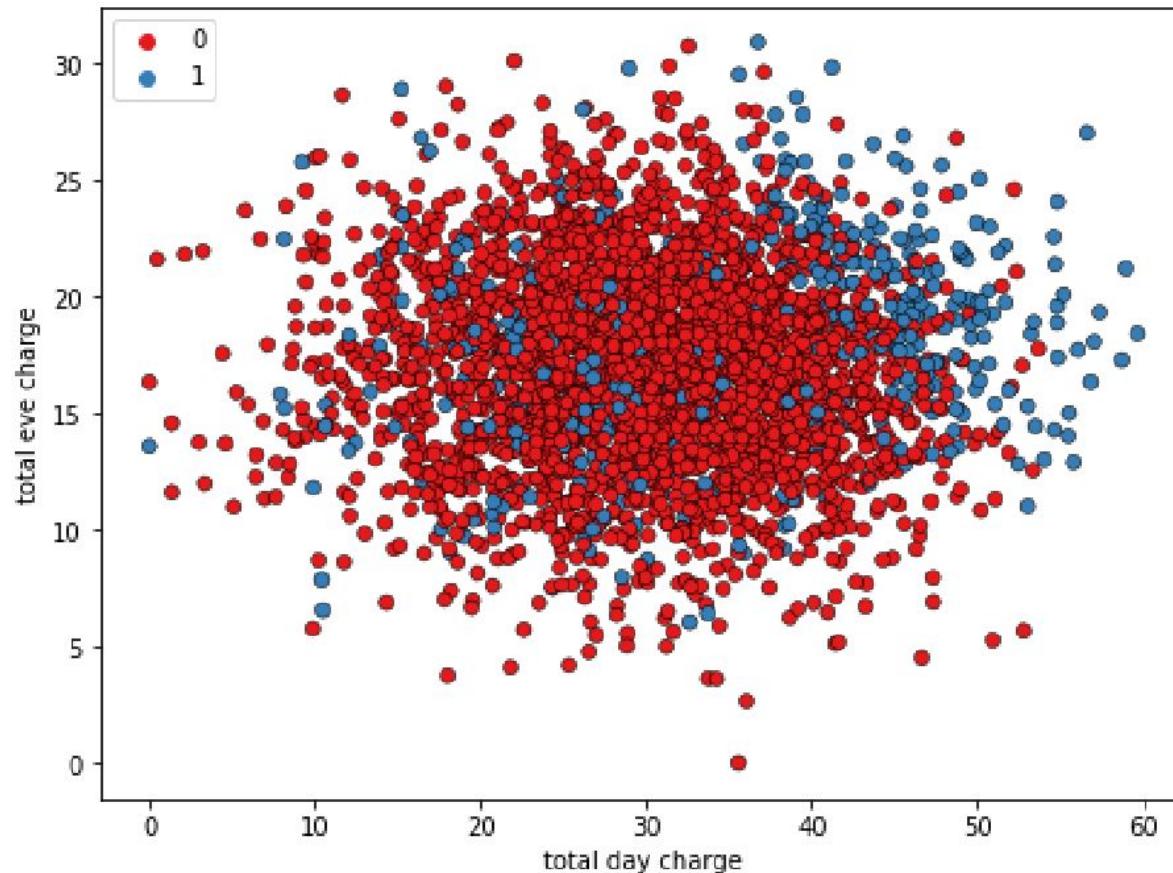


# k-Nearest Neighbors

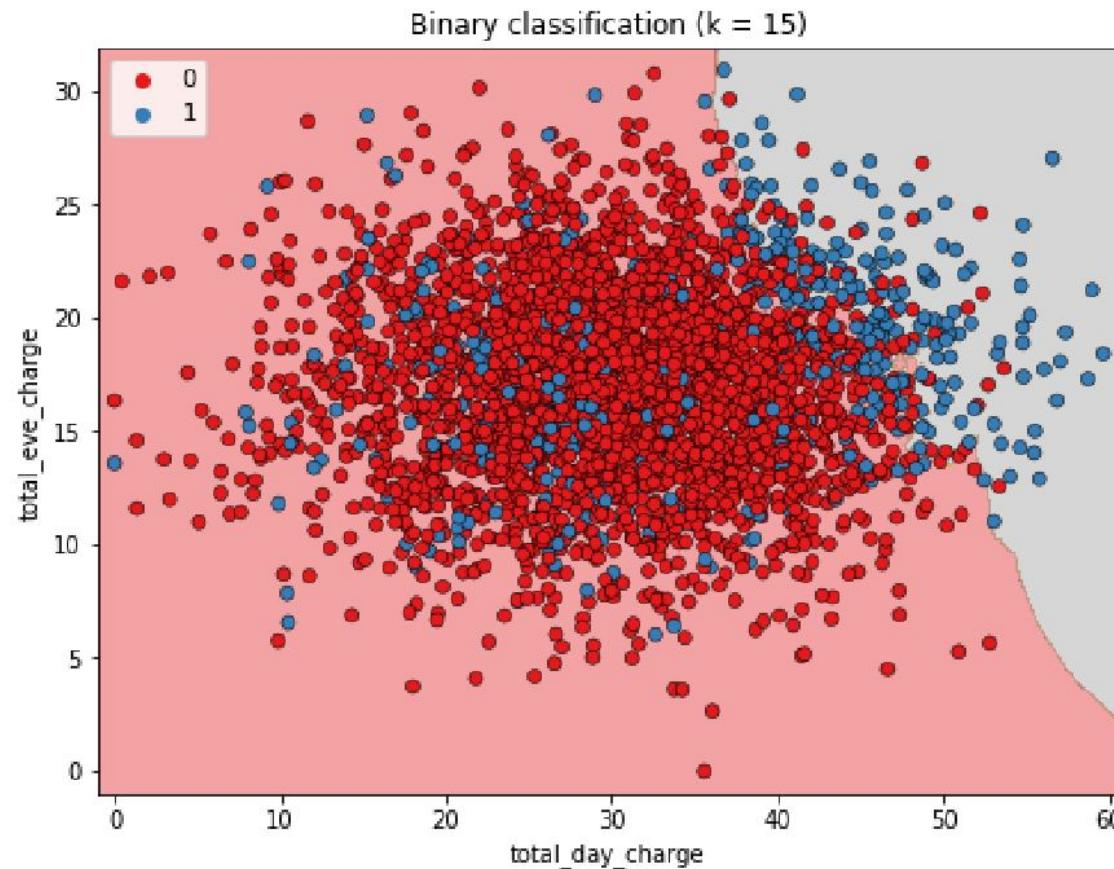
- K=5



# KNN Intuition



# KNN Intuition



# Using scikit-learn to fit a classifier

```
from sklearn.neighbors import KNeighborsClassifier  
X = churn_df[["total_day_charge", "total_eve_charge"]].values  
y = churn_df["churn"].values  
print(X.shape, y.shape)
```

```
(3333, 2), (3333,)
```

```
knn = KNeighborsClassifier(n_neighbors=15)  
knn.fit(X, y)
```

# Predicting on unlabeled data

```
X_new = np.array([[56.8, 17.5],  
                  [24.4, 24.1],  
                  [50.1, 10.9]])  
  
print(X_new.shape)
```

```
(3, 2)
```

```
predictions = knn.predict(X_new)  
print('Predictions: {}'.format(predictions))
```

```
Predictions: [1 0 0]
```

# Measuring Model Performance

# Measuring model performance

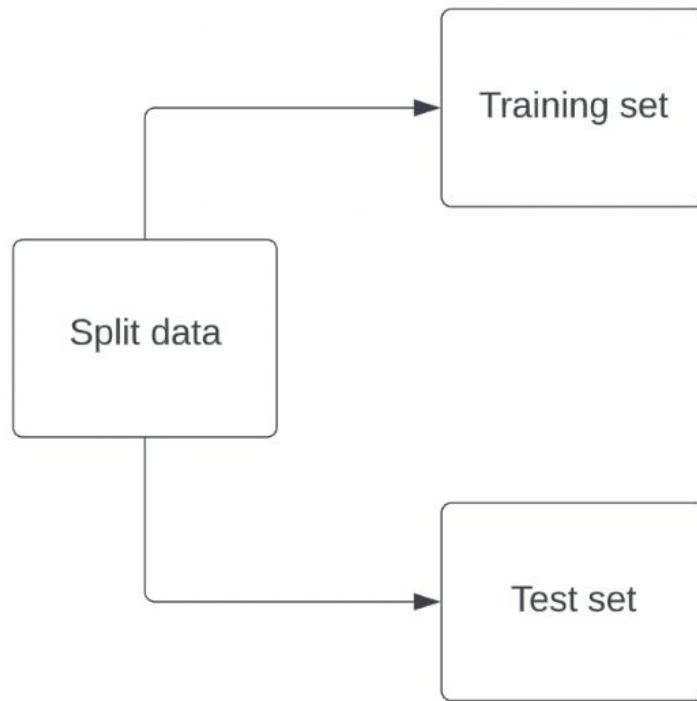
- In classification, accuracy is a commonly used metric
- Accuracy:

$$\frac{\text{correct predictions}}{\text{total observations}}$$

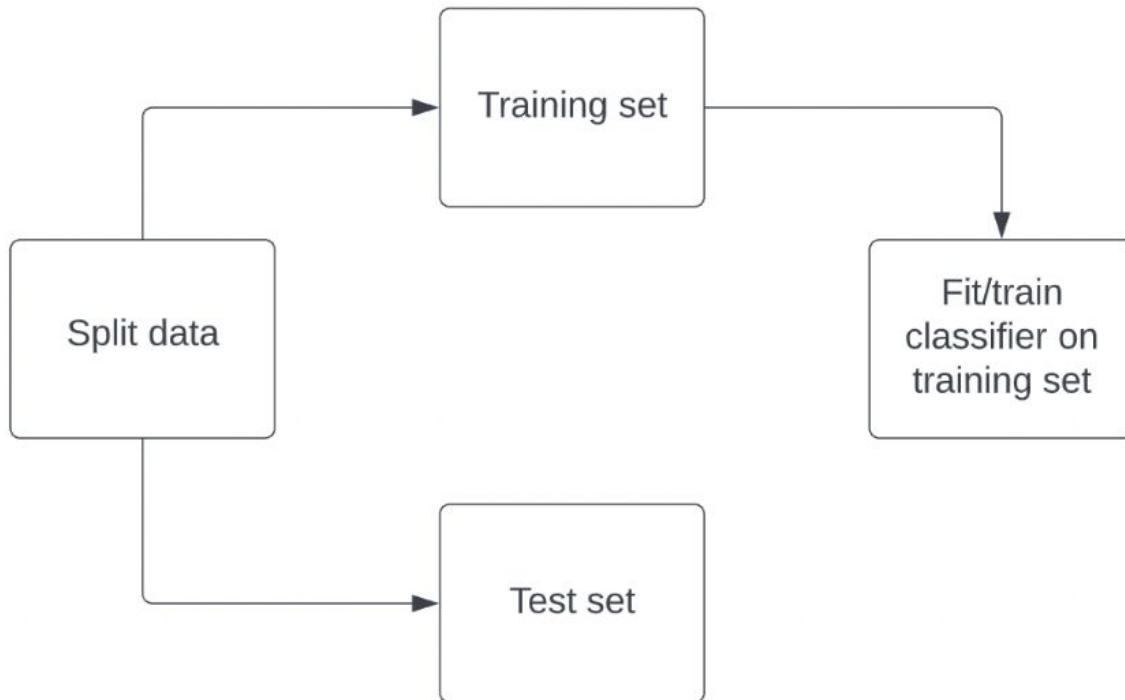
# Measuring model performance

- How do we measure accuracy?
- Could compute accuracy on the data used to fit the classifier
- NOT indicative of ability to generalize

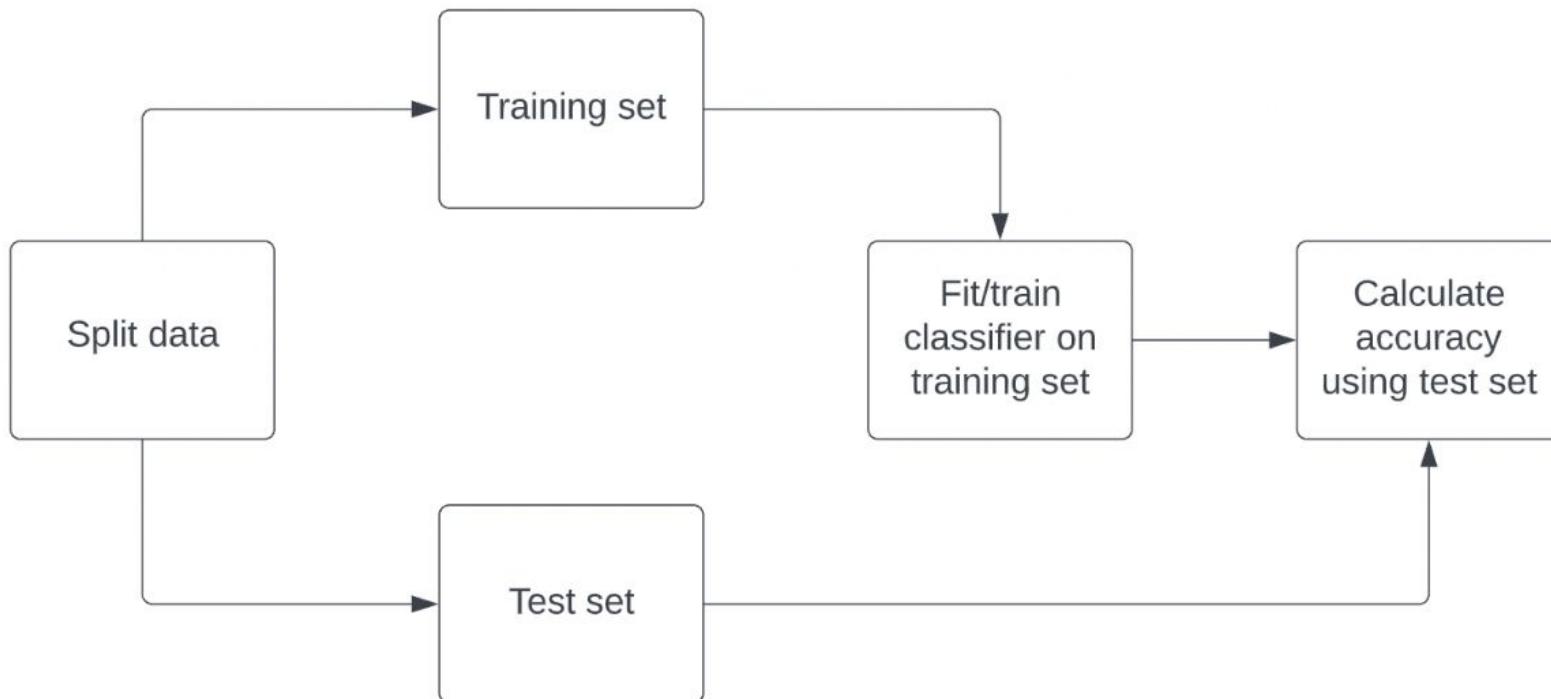
# Computing accuracy



# Computing accuracy



# Computing accuracy



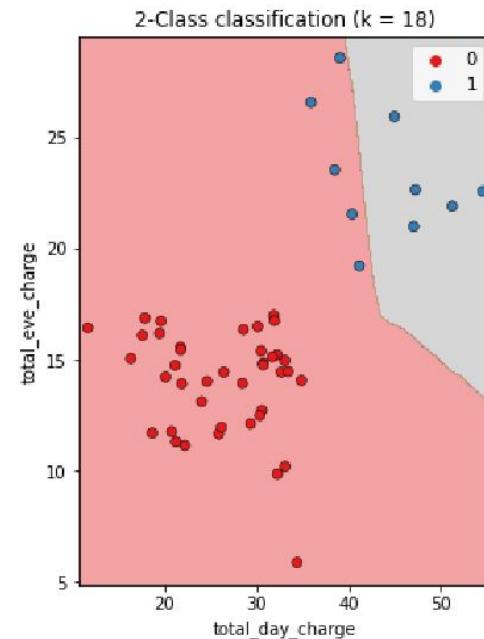
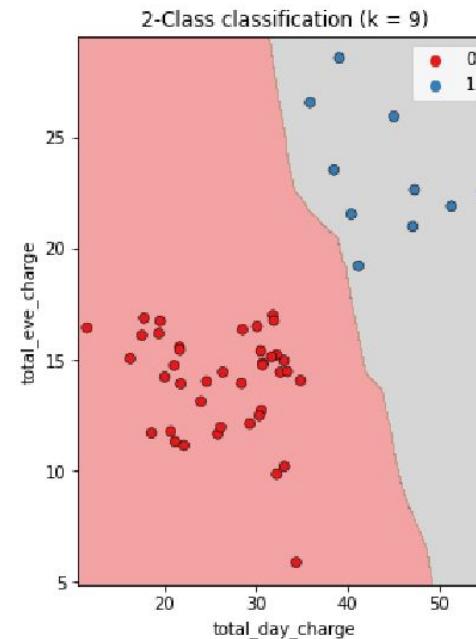
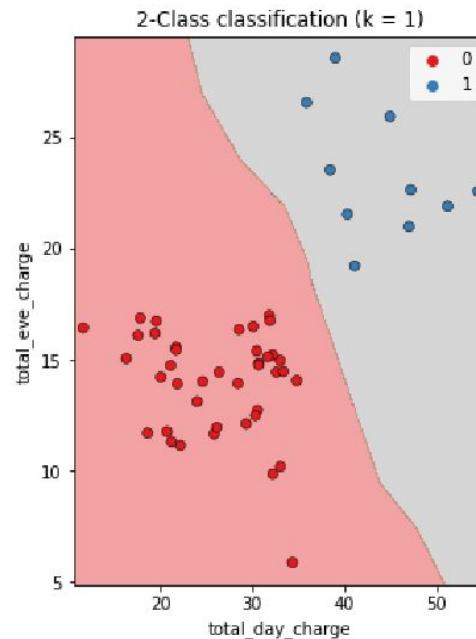
# Train/Test Split

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
                                                    random_state=21, stratify=y)
knn = KNeighborsClassifier(n_neighbors=6)
knn.fit(X_train, y_train)
print(knn.score(X_test, y_test))
```

0.8800599700149925

# Model Complexity

- Larger  $k$  = less complex model = can cause underfitting
- Smaller  $k$  = more complex model = can lead to overfitting



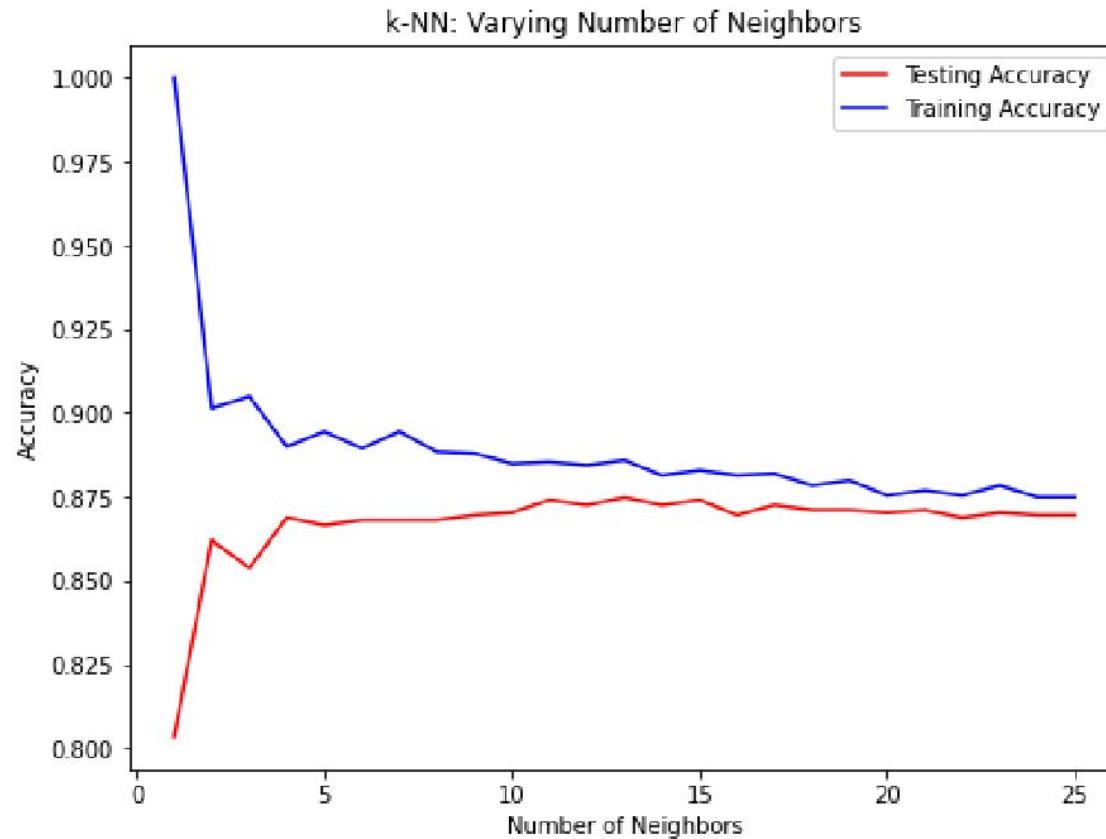
# Model complexity and over/underfitting

```
train_accuracies = []
test_accuracies = []
neighbors = np.arange(1, 26)
for neighbor in neighbors:
    knn = KNeighborsClassifier(n_neighbors=neighbor)
    knn.fit(X_train, y_train)
    train_accuracies[neighbor] = knn.score(X_train, y_train)
    test_accuracies[neighbor] = knn.score(X_test, y_test)
```

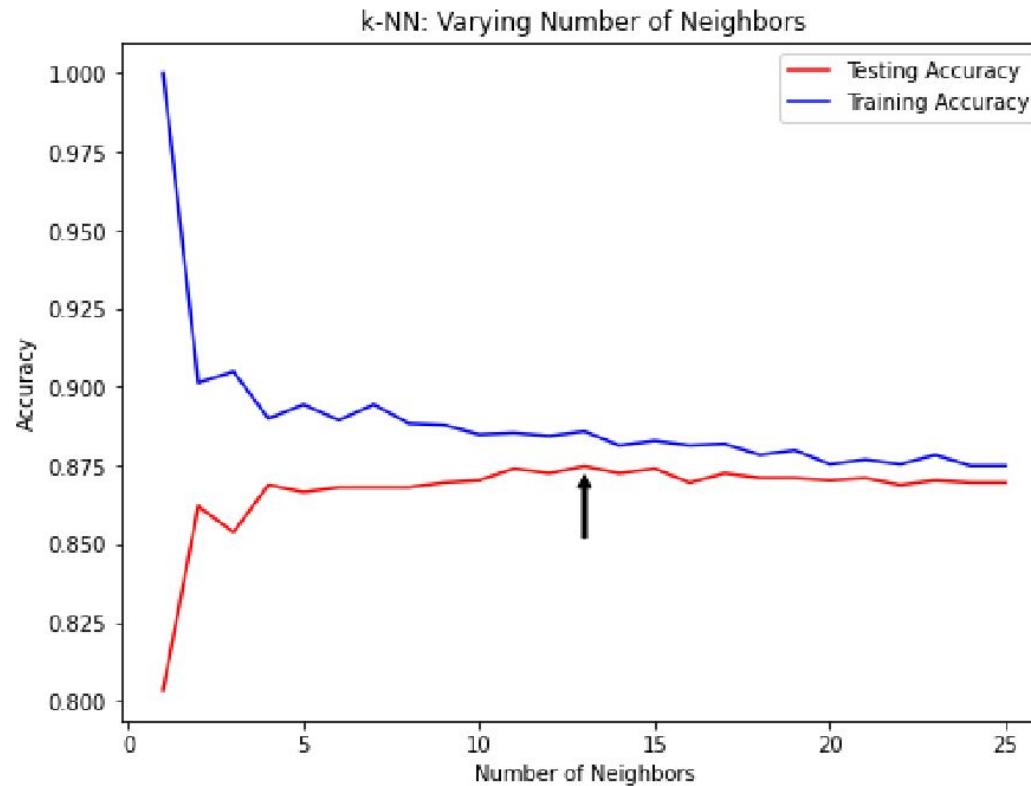
# Plotting our Results

```
plt.figure(figsize=(8, 6))
plt.title("KNN: Varying Number of Neighbors")
plt.plot(neighbors, train_accuracies.values(), label="Training Accuracy")
plt.plot(neighbors, test_accuracies.values(), label="Testing Accuracy")
plt.legend()
plt.xlabel("Number of Neighbors")
plt.ylabel("Accuracy")
plt.show()
```

# Model Complexity Curve



# Model Complexity Curve



*Practice*

## *Classification with scikit-learn*

*Classifying the churn status of a telecom company's customers.*

*The Classification Challenge*

# Regression

# Predicting blood glucose levels

```
import pandas as pd  
diabetes_df = pd.read_csv("diabetes.csv")  
print(diabetes_df.head())
```

	pregnancies	glucose	triceps	insulin	bmi	age	diabetes
0	6	148	35	0	33.6	50	1
1	1	85	29	0	26.6	31	0
2	8	183	0	0	23.3	32	1
3	1	89	23	94	28.1	21	0
4	0	137	35	168	43.1	33	1

# Creating feature and target arrays

```
X = diabetes_df.drop("glucose", axis=1).values  
y = diabetes_df["glucose"].values  
print(type(X), type(y))
```

```
<class 'numpy.ndarray'> <class 'numpy.ndarray'>
```

# Making predictions from a single feature

```
X_bmi = X[:, 3]  
print(y.shape, X_bmi.shape)
```

```
(752,) (752,)
```

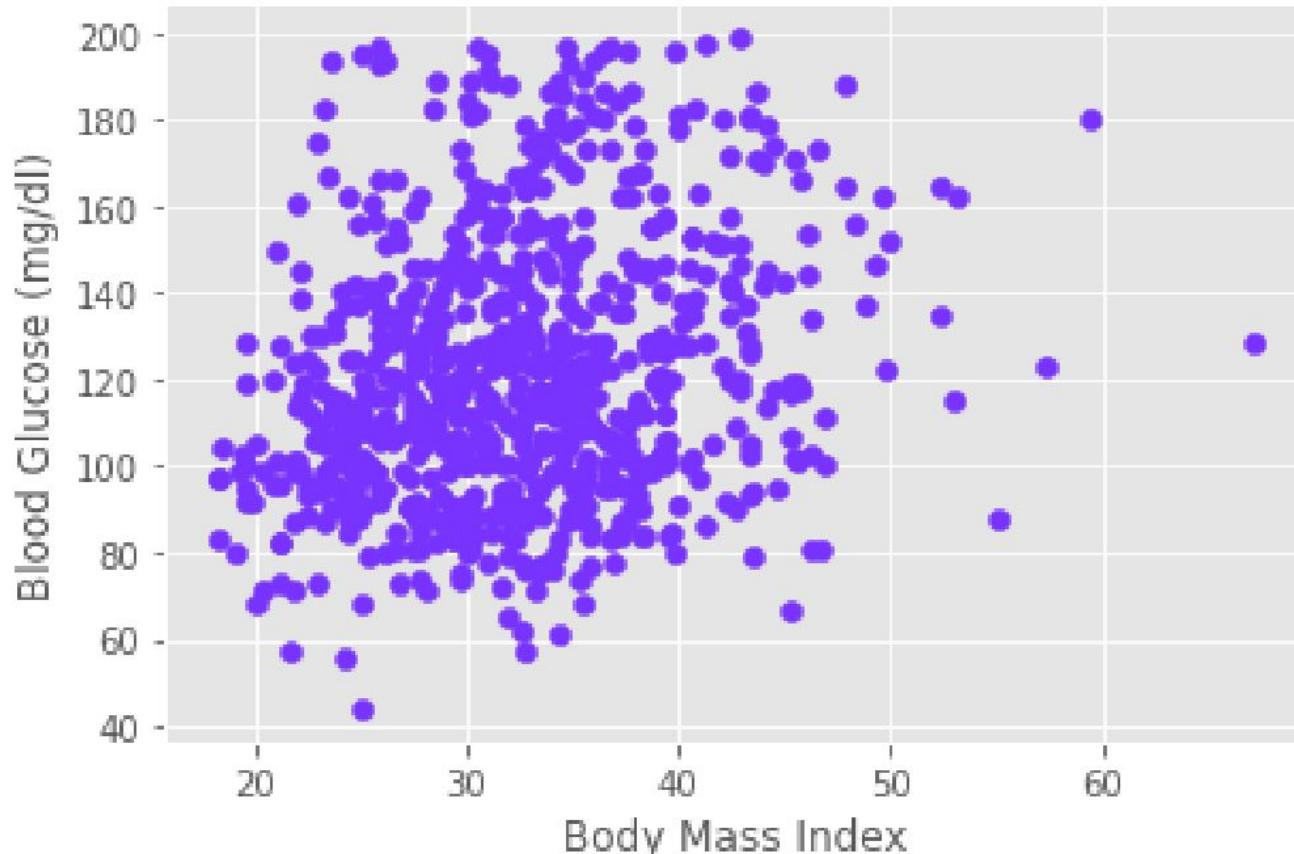
```
X_bmi = X_bmi.reshape(-1, 1)  
print(X_bmi.shape)
```

```
(752, 1)
```

# Plotting glucose vs. body mass index

```
import matplotlib.pyplot as plt
plt.scatter(X_bmi, y)
plt.ylabel("Blood Glucose (mg/dL)")
plt.xlabel("Body Mass Index")
plt.show()
```

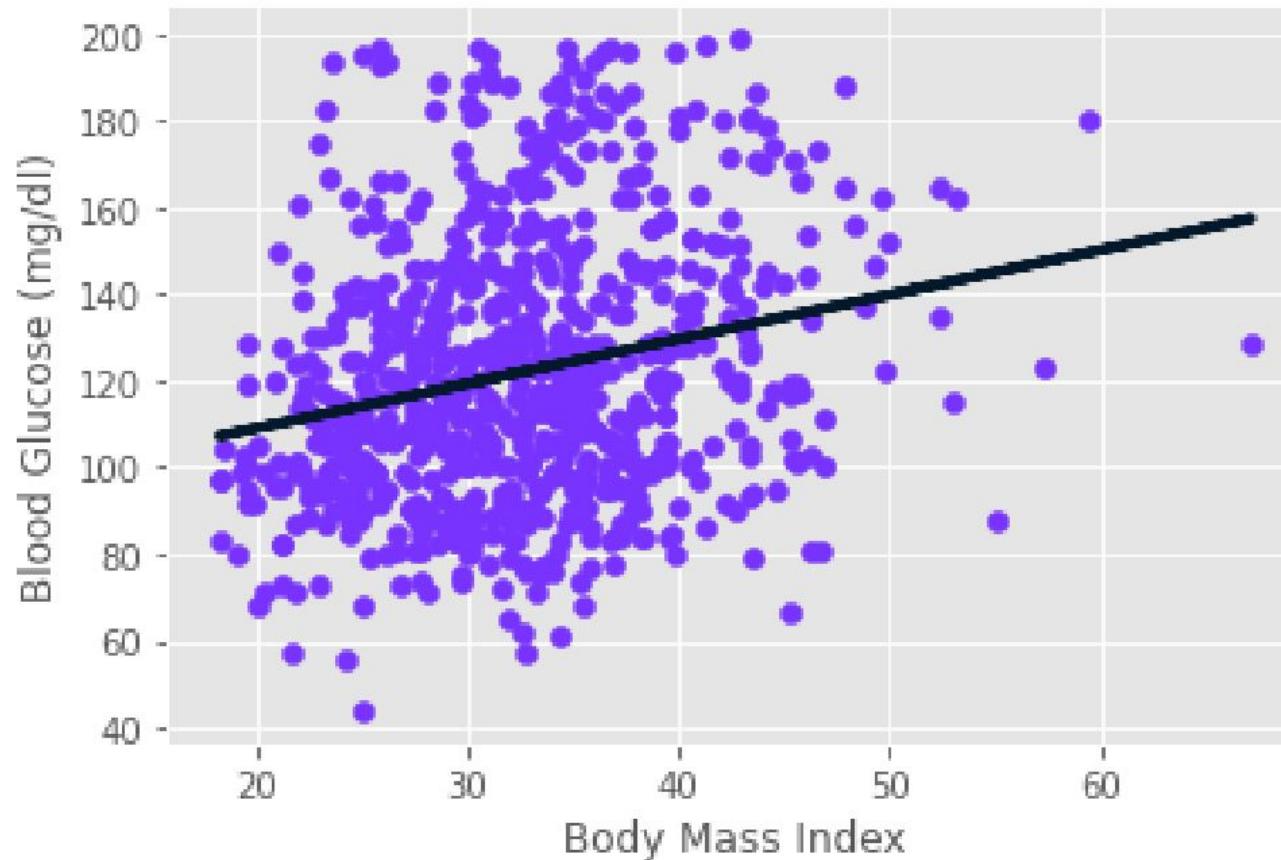
# Plotting glucose vs. body mass index



# Fitting a regression model

```
from sklearn.linear_model import LinearRegression  
reg = LinearRegression()  
reg.fit(X_bmi, y)  
predictions = reg.predict(X_bmi)  
plt.scatter(X_bmi, y)  
plt.plot(X_bmi, predictions)  
plt.ylabel("Blood Glucose (mg/dL)")  
plt.xlabel("Body Mass Index")  
plt.show()
```

# Fitting a regression model

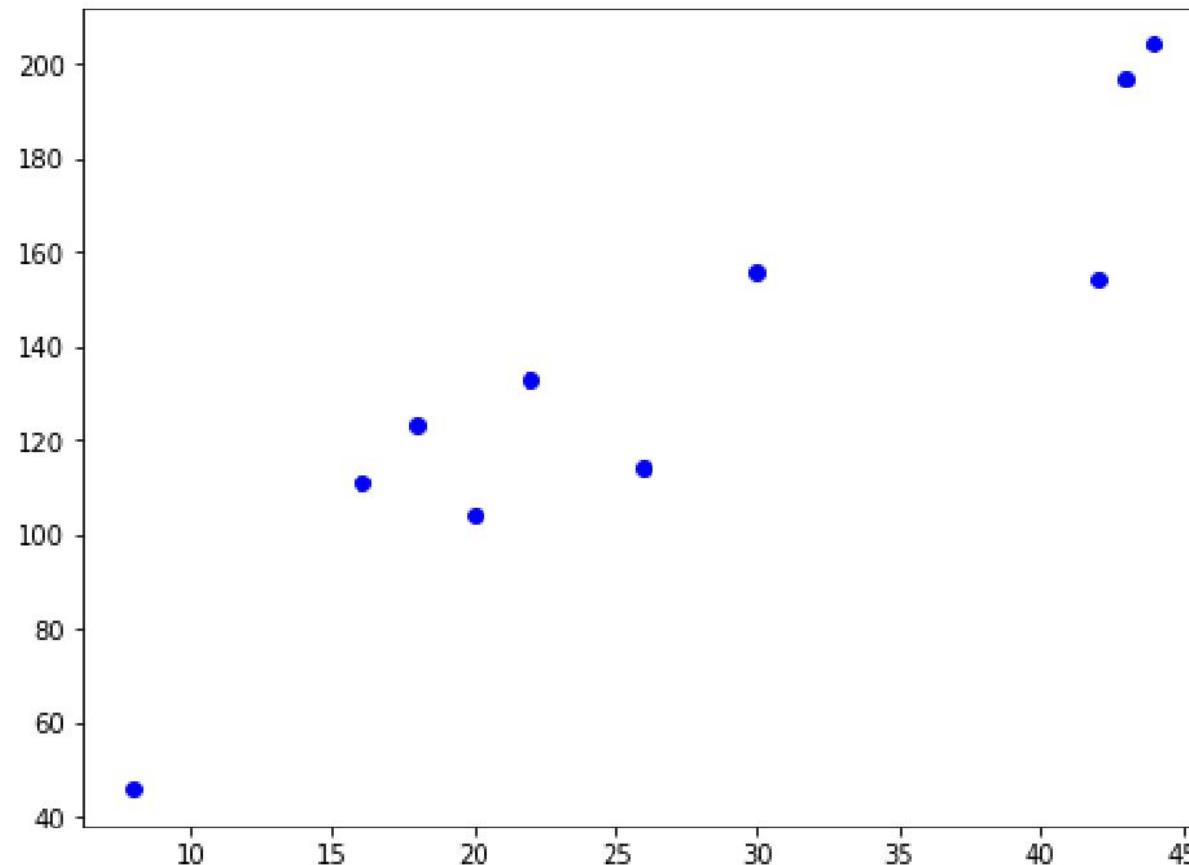


# The Basics of Linear Regression

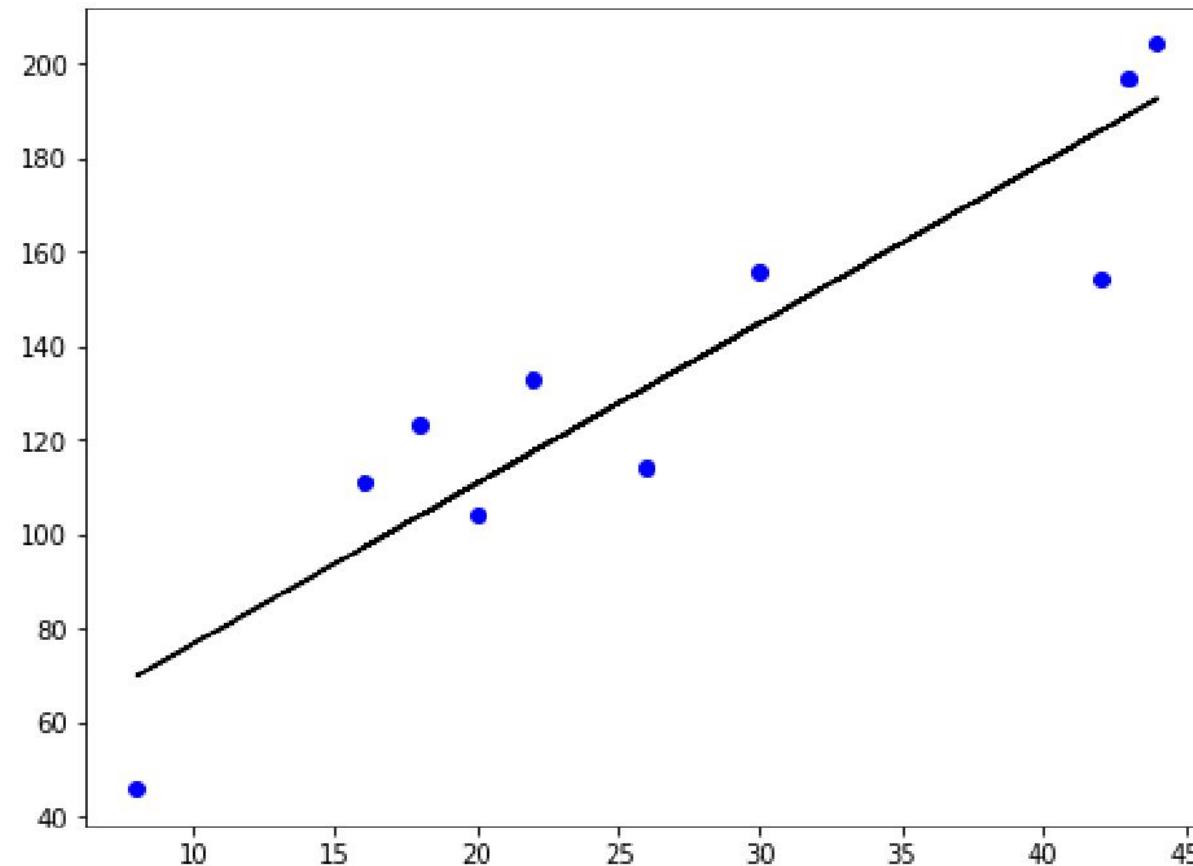
# Regression Mechanics

- $y = ax + b$ 
  - Simple linear regression uses one feature
    - $y$  = target
    - $x$  = single feature
    - $a, b$  = parameters/coefficients of the model - slope, intercept
- How do we choose  $a$  and  $b$ ?
  - Define an error function for any given line
  - Choose the line that minimizes the error function
- Error function = loss function = cost function

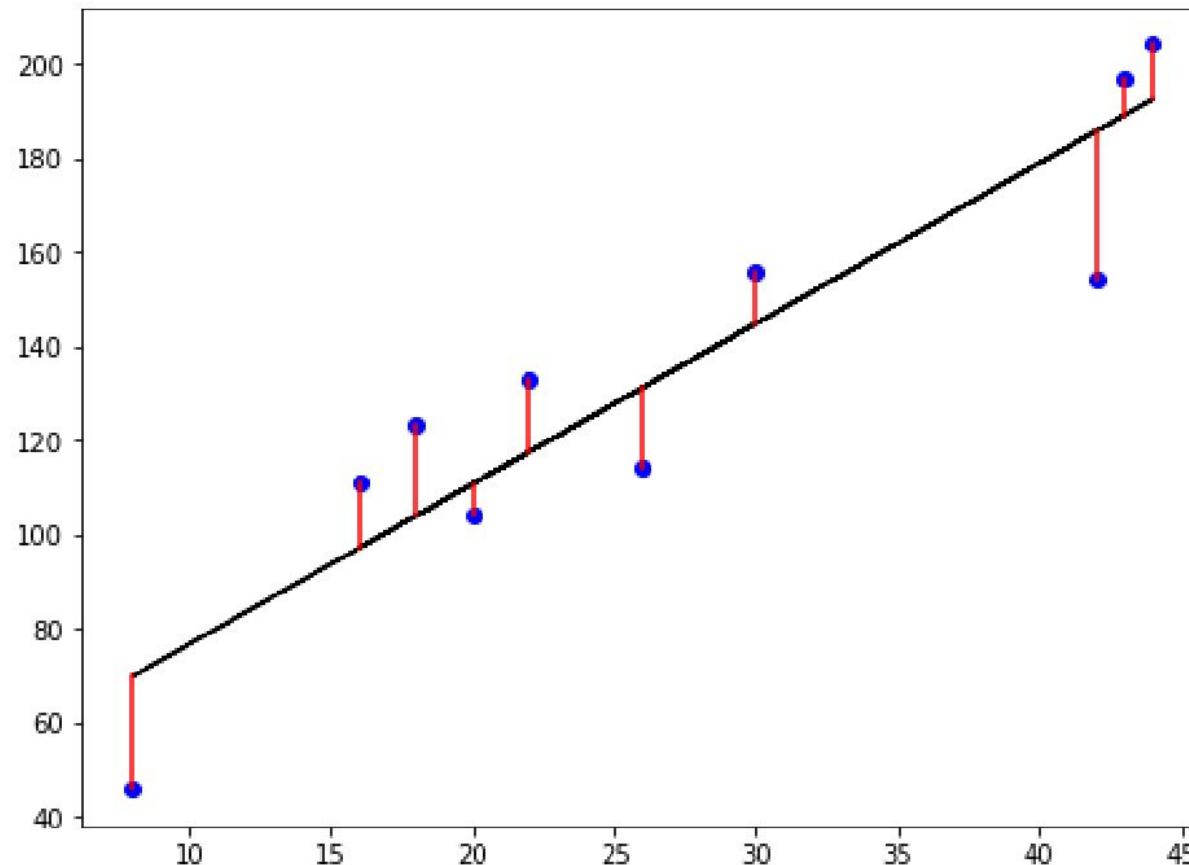
# The loss function



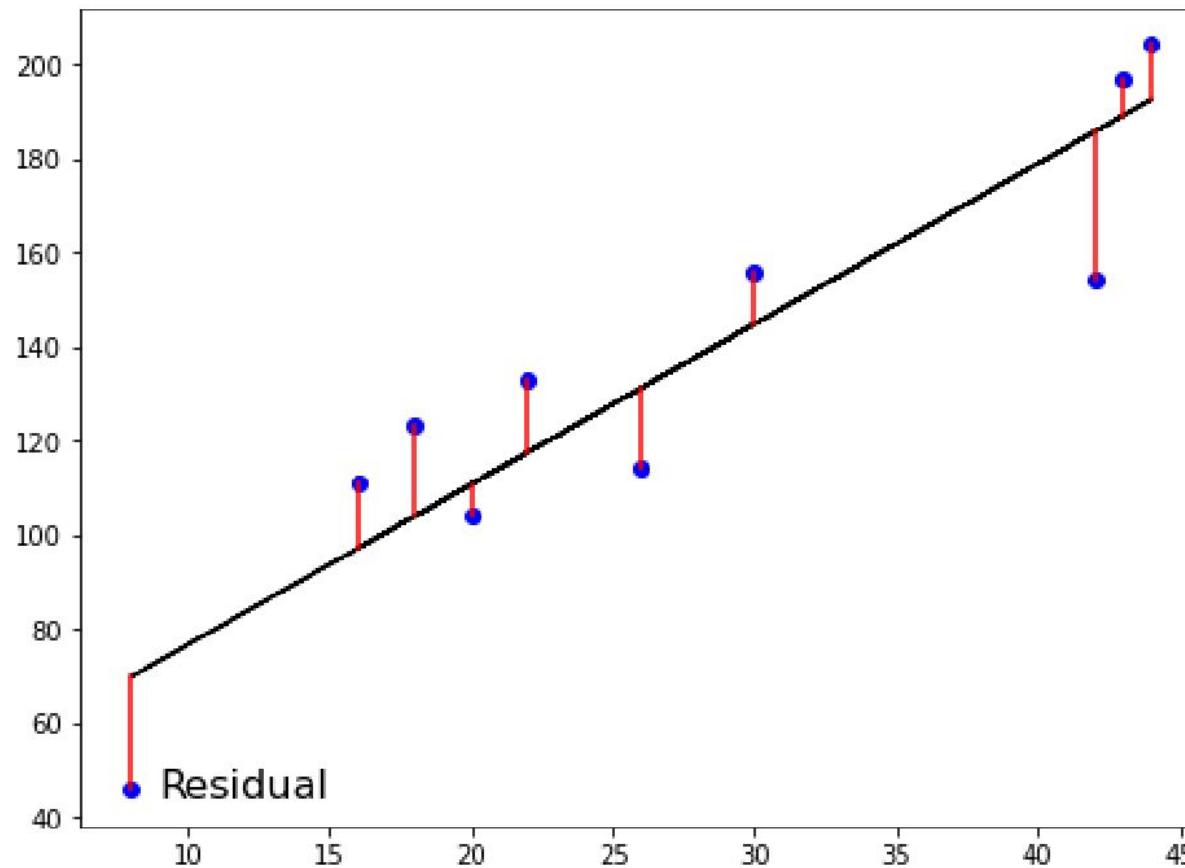
# The loss function



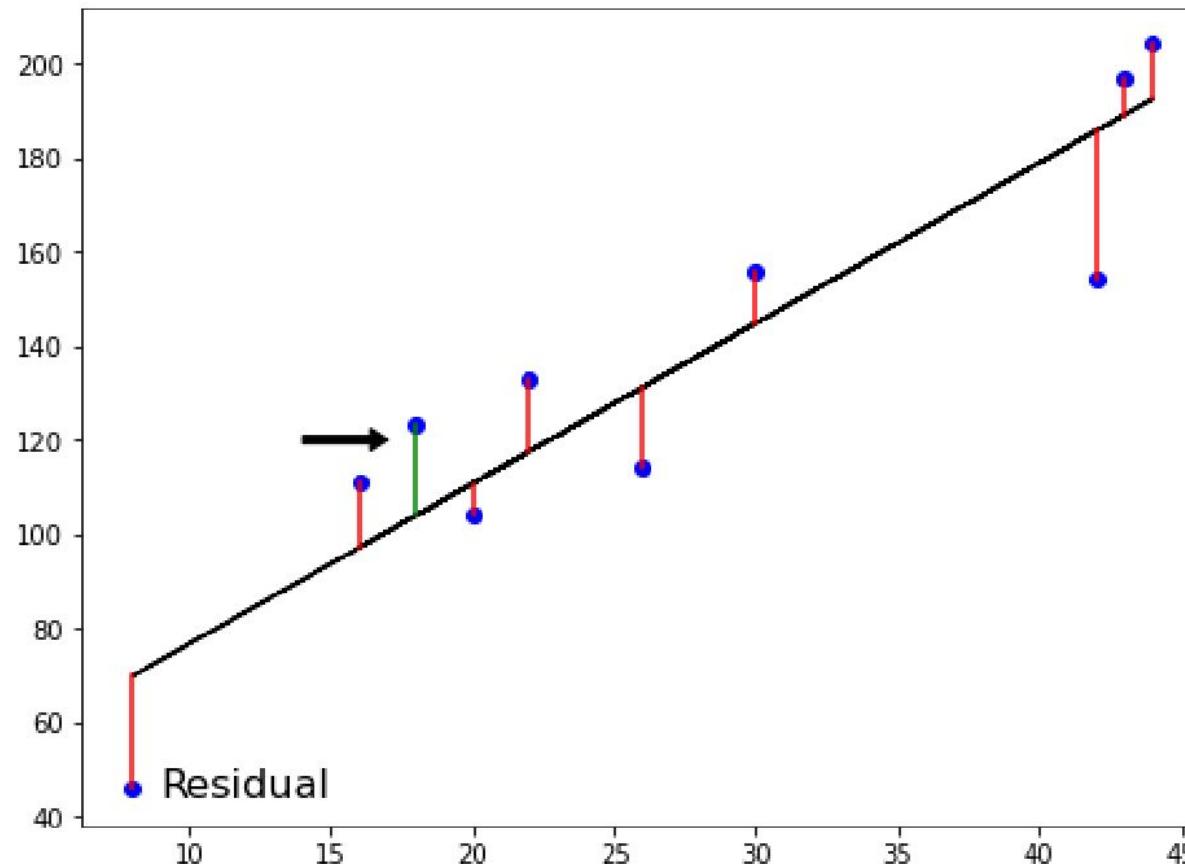
# The loss function



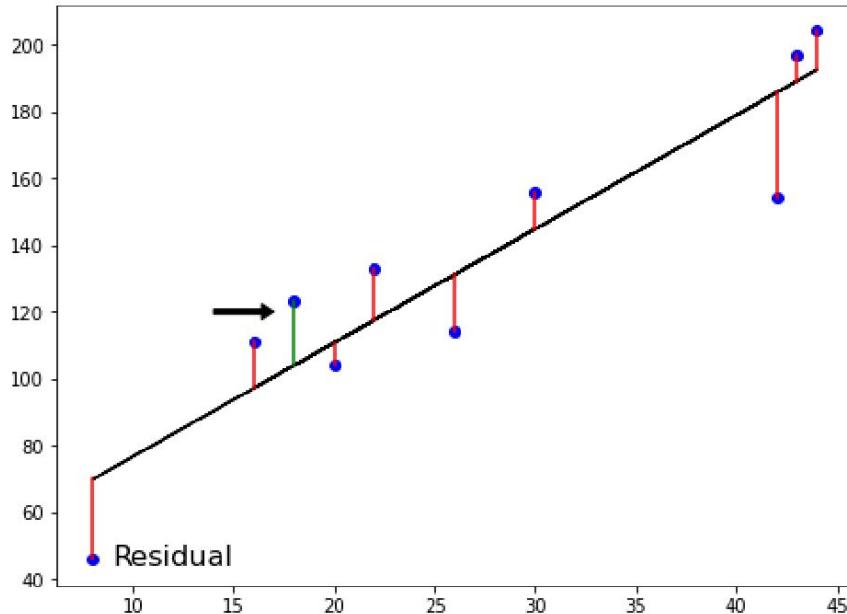
# The loss function



# The loss function



# Ordinary Least Squares



$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Ordinary Least Squares (OLS): minimize RSS

# Linear regression in higher dimensions

$$y = a_1x_1 + a_2x_2 + b$$

- To fit a linear regression model here:
  - Need to specify 3 variables:  $a_1, a_2, b$
- In higher dimensions:
  - Known as multiple regression
  - Must specify coefficients for each feature and the variable b

$$y = a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_nx_n + b$$

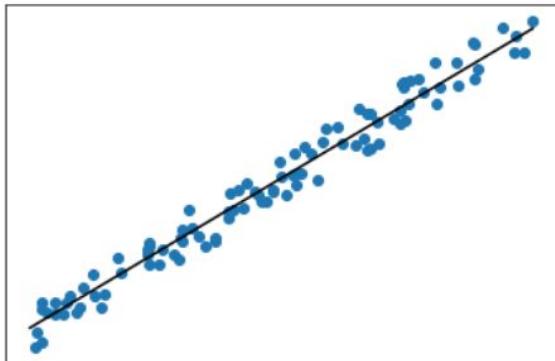
- scikit-learn works exactly the same way:
  - Pass two arrays: features and target

# Linear regression using all features

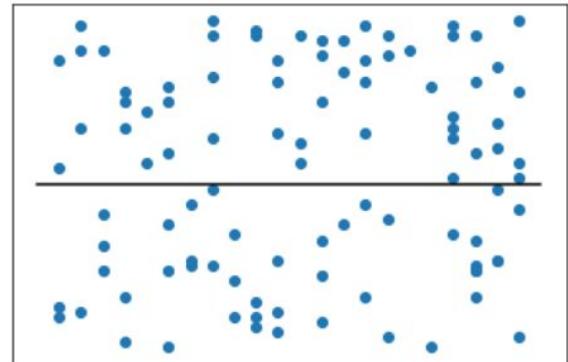
```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
                                                    random_state=42)
reg_all = LinearRegression()
reg_all.fit(X_train, y_train)
y_pred = reg_all.predict(X_test)
```

# R-squared

- $R^2$ : quantifies the variance in target values explained by the features
  - Values range from 0 to 1
- High  $R^2$ :



- Low  $R^2$ :



# R-squared in scikit-learn

```
reg_all.score(X_test, y_test)
```

```
0.356302876407827
```

# R-squared in scikit-learn

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- MSE is measured in target units, squared

$$RMSE = \sqrt{MSE}$$

- Measure RMSE in the same units at the target variable

# RMSE in scikit-learn

```
from sklearn.metrics import mean_squared_error  
mean_squared_error(y_test, y_pred, squared=False)
```

```
24.028109426907236
```

# Cross-validation

# Cross-validation Motivation

- Model performance is dependent on the way we split up the data
- Not representative of the model's ability to generalize to unseen data
- Solution: Cross-validation!

# Cross-validation Basics

Split 1

Fold 1

Fold 2

Fold 3

Fold 4

Fold 5

Metric 1

Training Data      Test Data

# Cross-validation Basics

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 1						
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2

Training Data      Test Data

# Cross-validation Basics

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 1						
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3

Training Data      Test Data

# Cross-validation Basics

Split 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3
Split 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 4

Training Data      Test Data

# Cross-validation Basics

Split 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3
Split 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 4
Split 5	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 5

Training Data      Test Data

# Cross-validation and Model Performance

- 5 folds = 5-fold CV
- 10 folds = 10-fold CV
- $k$  folds =  $k$ -fold CV
- More folds = More computationally expensive

# Cross-validation in scikit-learn

```
from sklearn.model_selection import cross_val_score, KFold
kf = KFold(n_splits=6, shuffle=True, random_state=42)
reg = LinearRegression()
cv_results = cross_val_score(reg, X, y, cv=kf)
```

# Evaluating Cross-validation Performance

```
print(cv_results)
```

```
[0.70262578, 0.7659624, 0.75188205, 0.76914482, 0.72551151, 0.73608277]
```

```
print(np.mean(cv_results), np.std(cv_results))
```

```
0.7418682216666667 0.023330243960652888
```

```
print(np.quantile(cv_results, [0.025, 0.975]))
```

```
array([0.7054865, 0.76874702])
```

# Regularized Regression

# Why Regularize

- Recall: Linear regression minimizes a loss function
- It chooses a coefficient, a, for each feature variable, plus b
- Large coefficients can lead to overfitting
- Regularization: Penalize large coefficients

# Why Regularize

- Loss function = OLS loss function +

$$\alpha * \sum_{i=1}^n a_i^2$$

- Ridge penalizes large positive or negative coefficients
- $\alpha$ : parameter we need to choose
- Picking  $\alpha$  is similar to picking  $k$  in KNN
- Hyperparameter: variable used to optimize model parameters
- $\alpha$  controls model complexity
  - $\alpha = 0$  = OLS (Can lead to overfitting)
  - Very high  $\alpha$ : Can lead to underfitting

# Ridge regression in scikit-learn

```
from sklearn.linear_model import Ridge
scores = []
for alpha in [0.1, 1.0, 10.0, 100.0, 1000.0]:
    ridge = Ridge(alpha=alpha)
    ridge.fit(X_train, y_train)
    y_pred = ridge.predict(X_test)
    scores.append(ridge.score(X_test, y_test))
print(scores)
```

```
[0.2828466623222221, 0.28320633574804777, 0.2853000732200006,
 0.26423984812668133, 0.19292424694100963]
```

# Lasso regression

Loss function = OLS loss function +

$$\alpha * \sum_{i=1}^n |a_i|$$

# Lasso regression in scikit-learn

```
from sklearn.linear_model import Lasso
scores = []
for alpha in [0.01, 1.0, 10.0, 20.0, 50.0]:
    lasso = Lasso(alpha=alpha)
    lasso.fit(X_train, y_train)
    lasso_pred = lasso.predict(X_test)
    scores.append(lasso.score(X_test, y_test))
print(scores)
```

```
[0.99991649071123, 0.99961700284223, 0.93882227671069, 0.74855318676232, -0.05741034640016]
```

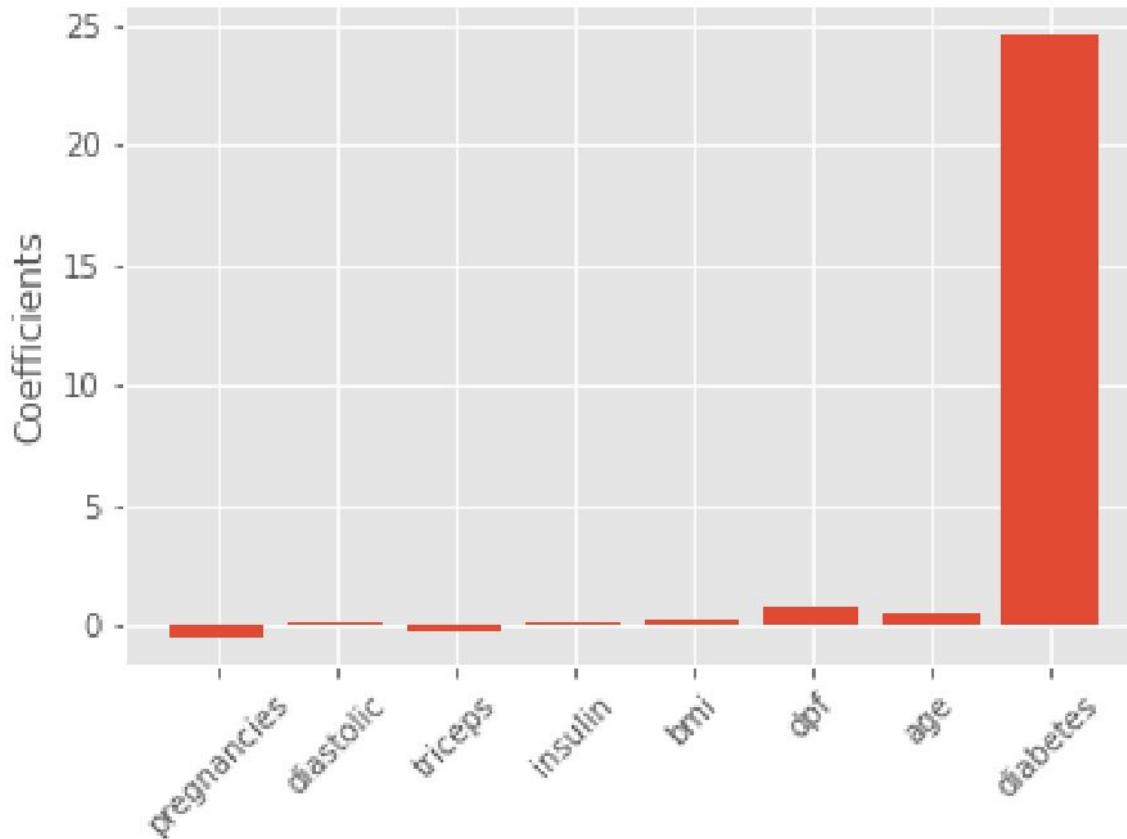
# Lasso regression for feature selection

- Lasso can select important features of a dataset
- Shrinks the coefficients of less important features to zero
- Features not shrunk to zero are selected by lasso

# Lasso for feature selection in scikit-learn

```
from sklearn.linear_model import Lasso
X = diabetes_df.drop("glucose", axis=1).values
y = diabetes_df["glucose"].values
names = diabetes_df.drop("glucose", axis=1).columns
lasso = Lasso(alpha=0.1)
lasso_coef = lasso.fit(X, y).coef_
plt.bar(names, lasso_coef)
plt.xticks(rotation=45)
plt.show()
```

# Lasso for feature selection in scikit-learn



*Practice*

## *Regression with scikit-learn*

*Predicting sales from advertising expenditure.*

*Introduction to Regression*

**Any  
questions?**

# THANKS

