

Homework #1

Course: *Machine Learning (CS405)* – Professor: *Qi Hao*

Due date: *11:59pm, September 23th, 2020*

NAM^E: 车^W锐
NUMBER: 12032207

Question 1

Consider the polynomial function

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Calculate the coefficients $\mathbf{w} = \{w_i\}$ that minimize its sum-of-squares error function. Here a suffix i or j denotes the index of a component, whereas $(x)^i$ denotes x raised to the power of i .

Question 2

Suppose that we have three colored boxes r (red), b (blue), and g (green). Box r contains 3 apples, 4 oranges, and 3 limes, box b contains 1 apple, 1 orange, and 0 limes, and box g contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities $p(r) = 0.2$, $p(b) = 0.2$, $p(g) = 0.6$, and a piece of fruit is removed from the box (with equal probability of selecting any of the items in the box), then what is the probability of selecting an apple? If we observe that the selected fruit is in fact an orange, what is the probability that it came from the green box?

Question 3

Given two statistically independent variables x and z , show that the mean and variance of their sum satisfies

$$\mathbb{E}[x + z] = \mathbb{E}[x] + \mathbb{E}[z]$$

$$\text{var}[x + z] = \text{var}[x] + \text{var}[z]$$

Question 1 :

Assume that: M means the order of the polynomial

N means N samples

t_n means the real value of the input of $x^{(n)}$

$y(x, w)$ means the predict value of the input of $x^{(n)}$

$$\text{Error Function: } E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2$$

$$\text{Gradient Descent: } w_j := w_j - \alpha \frac{\partial}{\partial w_j} E(w)$$

$$\frac{\partial}{\partial w_j} E(w) = \sum_{n=1}^N (y(x_n, w) - t_n) \cdot x_j^{(n)}$$

Therefore, the w_j that minimize its sum-of-squares-error function is:

$$w_j := w_j - \alpha \cdot \sum_{n=1}^N (y(x_n, w) - t_n) \cdot x_j^{(n)}, j \in M$$

Question 2:

Assume event F means the probability of selecting fruit from box.

a is apple, o is orange, l is limes

Assume event B means the probability of selecting box

$$\begin{aligned}1. \quad p(F=a) &= \sum_B p(F=a|B) \\&= p(F=a|B=r)p(B=r) + p(F=a|B=b)p(B=b) \\&\quad + p(F=a|B=g)p(B=g) \\&= \frac{2}{10} \cdot \frac{3}{10} + \frac{2}{10} \cdot \frac{1}{2} + \frac{6}{10} \cdot \frac{3}{10} \\&= \frac{34}{100} = 0.34\end{aligned}$$

2. According to Bayes' theorem:

$$p(B=g|F=o) = \frac{p(F=o|B=g)p(B=g)}{p(F=o)}$$

$$p(F=o|B=g) = \frac{3}{10} = 0.3$$

$$p(B=g) = 0.6$$

$$\begin{aligned}p(F=o) &= \sum_B p(F=o|B) \\&= p(F=o|B=r)p(B=r) + p(F=o|B=b)p(B=b) \\&\quad + p(F=o|B=g)p(B=g) \\&= \frac{4}{10} \times \frac{2}{10} + \frac{1}{2} \times \frac{2}{10} + \frac{3}{10} \times \frac{6}{10} \\&= \frac{36}{100} = 0.36\end{aligned}$$

$$\begin{aligned}p(B=g|F=o) &= \frac{p(F=o|B=g)p(B=g)}{p(F=o)} \\&= \frac{0.3 \times 0.6}{0.36} = \frac{1}{2}\end{aligned}$$

Question 3 :

$$E(x+z) = E(x) + E(z)$$

Discrete :

$$E[x] = \sum_{i=1}^M x_i p(x_i)$$

$$E[z] = \sum_{j=1}^N z_j P(z_j)$$

$$E[x+z] = \sum_{i=1}^M \sum_{j=1}^N (x_i + z_j) P(x_i z_j)$$

$$= \sum_{i=1}^M \sum_{j=1}^N (x_i + z_j) P(x_i) P(z_j)$$

$$= \sum_{i=1}^M x_i p(x_i) + \sum_{j=1}^N z_j P(z_j)$$

$$= E[x] + E[z]$$

Continuous :

$$E[x] = \int x f(x) dx$$

$$E[z] = \int z f(z) dz$$

$$E[x+z] = \iint (x+z) f(xz) dx dz$$

$$= \iint (x+z) f(x) f(z) dx dz$$

$$= \iint x f(x) f(z) dx dz + \iint z f(x) f(z) dx dz$$

$$= \int x f(x) dx + \int z f(z) dz$$

$$= E[x] + E[z]$$

$$\text{var}[x+z] = \text{var}[x] + \text{var}[z]$$

$$\text{var}[x] = E(x^2) - (Ex)^2$$

$$\text{var}[z] = E(z^2) - (Ez)^2$$

$$\begin{aligned}\text{var}[x+z] &= E((x+z)^2) - (E(x+z))^2 \\&= E(x^2 + 2xz + z^2) - [E(x)^2 + 2ExEz + (Ez)^2] \\&= Ex^2 + 2ExEz + Ez^2 - (Ex)^2 - 2ExEz - (Ez)^2 \\&= Ex^2 - (Ex)^2 + Ez^2 - (Ez)^2 \\&= \text{var}[x] + \text{var}[z]\end{aligned}$$

Question 4

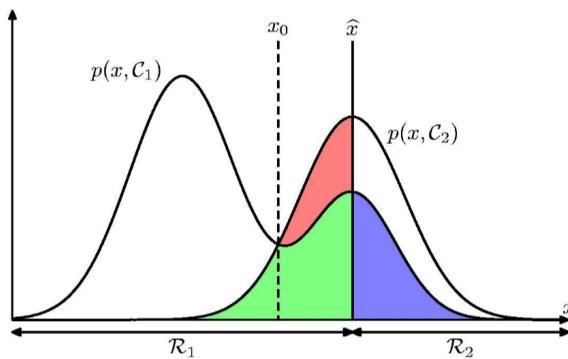
In probability theory and statistics, the Poisson distribution, is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time or space if these events occur with a known constant rate and independently of the time since the last event. If X is Poisson distributed, i.e. $X \sim \text{Poisson}(\lambda)$, its probability mass function takes the following form:

$$P(X|\lambda) = \frac{\lambda^X e^{-\lambda}}{X!}$$

It can be shown that if $\mathbb{E}(X) = \lambda$. Assume now we have n data points from $\text{Poisson}(\lambda) : \mathcal{D} = \{X_1, X_2, \dots, X_n\}$. Show that the sample mean $\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n X_i$ is the maximum likelihood estimate(MLE) of λ . If X is exponential distribution and its distribution density function is $f(x) = \frac{1}{\lambda} e^{-\frac{x}{\lambda}}$ for $x > 0$ and $f(x) = 0$ for $x \leq 0$. Show that the sample mean $\hat{\lambda} \frac{1}{n} \sum_{i=1}^n X_i$ is the maximum likelihood estimate(MLE) of λ .

Question 5

- (a) Write down the probability of classifying correctly $p(\text{correct})$ and the probability of misclassification $p(\text{mistake})$ according to the following chart.



- (b) For multiple target variables described by vector \mathbf{t} , the expected squared loss function is given by

$$\mathbb{E}[L(\mathbf{t}, \mathbf{y}(\mathbf{x}))] = \int \int \|\mathbf{y}(\mathbf{x}) - \mathbf{t}\|^2 p(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t}$$

Show that the function $\mathbf{y}(\mathbf{x})$ for which this expected loss is minimized given by $\mathbf{y}(\mathbf{x}) = \mathbb{E}_{\mathbf{t}}[\mathbf{t}|\mathbf{x}]$.

Hints. For a single target variable t , the loss is given by

$$\mathbb{E}[L] = \int \int \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

The result is as follows

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) dt = \mathbb{E}_t[t|\mathbf{x}]$$

Question 4 :

$$P(X|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$L(\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$$

$$= \frac{1}{x_1! \cdot x_2! \cdots x_n!} \cdot \lambda^{\sum_{i=1}^n x_i} e^{-\lambda n}$$

$$\ln L(\lambda) = \ln \frac{1}{x_1! \cdots x_n!} + \sum_{i=1}^n x_i \ln \lambda - \lambda n$$

$$\frac{d \ln L(\lambda)}{d \lambda} = \sum_{i=1}^n x_i \frac{1}{\lambda} - n = 0$$

$$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n x_i$$

Question 5 :

$$(a) P(\text{correct}) = \sum_{k=1}^K P(x \in R_k, C_k)$$

$$= \sum_{k=1}^K \int_{R_k} P(x, C_k) dx$$

$$P(\text{mistake}) = P(x \in R_1, C_2) + P(x \in R_2, C_1)$$

$$= \int_{R_1} P(x, C_2) dx + \int_{R_2} P(x, C_1) dx$$

$$(b) \quad E[L] = \iint (y(x) - t)^2 p(x, t) dx dt$$

$$E[L] = \iint [y^2(x) - 2y(x)t + t^2] p(x, t) dx dt$$

$$= y^2(x) \iint p(x, t) dx dt -$$

$$2y(x) \iint t p(x, t) dx dt +$$

$$\iint t^2 p(x, t) dx dt$$

$$= y^2(x) - 2y(x) \int t p(t|x) dt$$

$$\frac{\partial E[L]}{\partial y(x)} = 2y(x) - 2 \int t p(t|x) dt = 0$$

$$y(x) = \int t p(t|x) dt = E_t(t|x)$$

Because multiple target variables vector t is defined by $t[i]$, therefore vector $y(x) = E_t(t|x)$

Question 6

- (a) We defined the entropy based on a discrete random variable \mathbf{X} as

$$H[\mathbf{X}] = - \sum_i p(x_i) \ln p(x_i) \quad \text{the probability.}$$

Now consider the case that \mathbf{X} is a continuous random variable with the probability density function $p(x)$. The entropy is defined as

$$H[\mathbf{X}] = - \int p(x) \ln p(x) dx$$

Assume that \mathbf{X} follows Gaussian distribution with the mean μ and variance σ^2 , i.e.

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \text{Gaussian.}$$

Please derive its entropy $H[\mathbf{X}]$.

- (b) Write down the mutual information $I(\mathbf{y}|\mathbf{x})$. Then show the following equation

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

Question 6:

$$\begin{aligned}
 a) H[x] &= - \int p(x) \ln p(x) dx \\
 &= \int \frac{1}{\sqrt{2\lambda}\delta} e^{-\frac{(x-\mu)^2}{2\delta^2}} \cdot [\ln \sqrt{2\lambda}\delta + \frac{(x-\mu)^2}{2\delta^2}] dx \\
 &= \frac{\ln \sqrt{2\lambda}\delta}{\sqrt{2\pi}\delta} \int e^{-\frac{(x-\mu)^2}{2\delta^2}} dx + \frac{1}{\sqrt{2\lambda}\delta} \int e^{-\frac{(x-\mu)^2}{2\delta^2}} \cdot \frac{(x-\mu)^2}{2\delta^2} dx \\
 t &= \frac{x-\mu}{\delta}, \\
 &= \frac{\tau_2}{2} \frac{\ln 2\tau_1}{\lambda} + \frac{1}{\sqrt{2\lambda}} \underbrace{\int e^{-\frac{t^2}{2}} \frac{t^2}{2} dt}_{e^{-\frac{t^2}{2}} \cdot t^2 + \int e^{-\frac{t^2}{2}} dt = e^{-\frac{t^2}{2}} \cdot t^2 + 1} \\
 &= \frac{\tau_2}{2} \cdot \frac{\ln 2\lambda}{\lambda} + \frac{1}{\sqrt{2\lambda}} e^{-\frac{(x-\mu)^2}{2\delta^2}} \cdot \frac{(x-\mu)^2}{\delta^2} + \frac{1}{\sqrt{2\lambda}}
 \end{aligned}$$

$$b) I(x,y) = \sum_x \sum_y P(x,y) \ln \frac{P(x,y)}{P(x)P(y)}$$

$$H(x) = - \sum_x P(x) \ln P(x)$$

$$H(y) = - \sum_y P(y) \ln P(y)$$

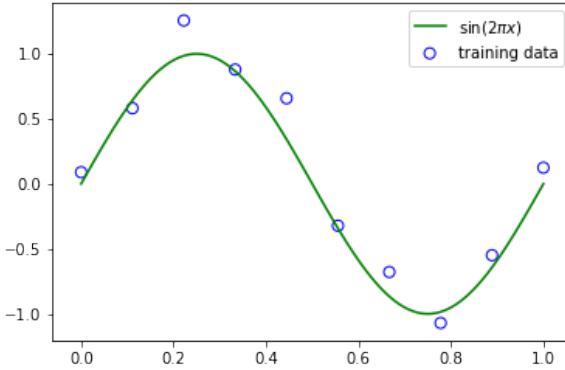
$$\begin{aligned}
 I(x,y) &= \sum_x \sum_y P(x,y) \ln \frac{P(x,y)}{P(y)} - \sum_x P(x) \ln P(x) \\
 &= \sum_x \sum_y P(y) P(x|y) \ln P(x|y) + H(x) \\
 &= -H(x|y) + H(x) = H(x) - H(x|y)
 \end{aligned}$$

$$\text{the same : } I(x,y) = H(y) - H(y|x)$$

Program Question

You should download the HW1_programQuestion.ipynb file first.

- (a) Plot the graph with given code, the result should be same as this.



- (b) On the basis of the results, you should try 0^{th} order polynomial, 1^{st} order polynomial, 3^{rd} order polynomial and some other order polynomial, show the results include fitting and over-fitting.
- (c) Plot the graph of the root-mean-square error.
- (d) Plot the graph of the predictive distribution resulting from a Bayesian treatment of polynomial curve fitting using an $M=9$ polynomial, with the fixed parameters $\alpha = 5 \times 10^{-3}$ and $\beta = 11.1$ (corresponding to the known noise variance).
- (e) Change the `sample_size` to 2, 3 or 10 times than before, explain the change of M .

Hints. You should install `matplotlib.pyplot`, and read classes `PolynomialFeature`, `LinearRegression`, and `BayesianRegression` in the file.

CS405 Machine Learning: HW 1 Preliminary

Name: 车凯威

ID: 12032207

0 Prepare the data

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [2]: def create_toy_data(func, sample_size, std):
    x = np.linspace(0, 1, sample_size)
    t = func(x) + np.random.normal(scale=std, size=x.shape)
    return x, t

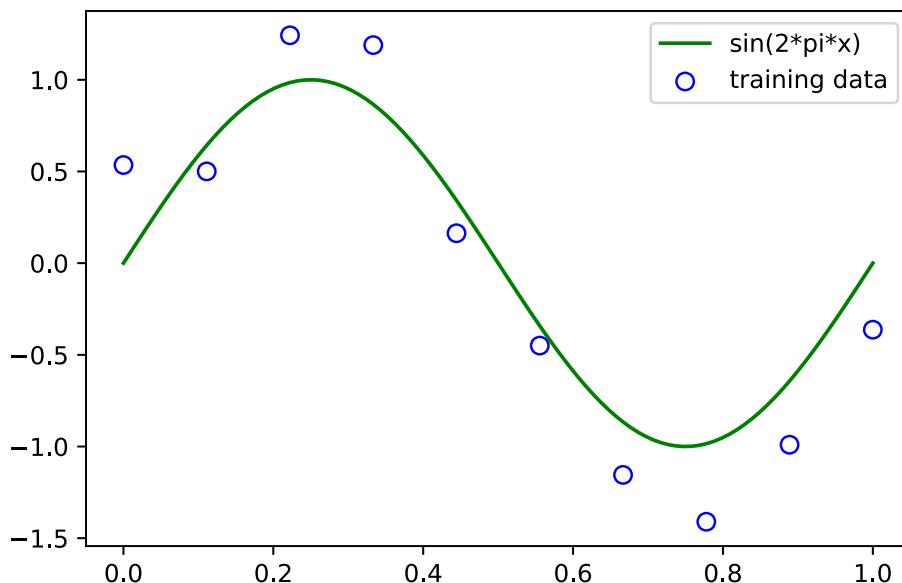
def func(x):
    return np.sin(2 * np.pi * x)

x_train, y_train = create_toy_data(func, 10, 0.25)
x_test = np.linspace(0, 1, 100)
y_test = func(x_test)
```

(a) Plot the graph with given code, the result should be same as this. `x_train` and `y_train` are the datas you need to create, `sample_size` is 10 and `std` is 0.25.

```
In [3]: # my code

plt.scatter(x_train, y_train, facecolor="none", edgecolor="b", s=50, label="training data")
plt.plot(x_test, y_test, c='g', label="sin(2*pi*x)")
plt.legend()
plt.show()
```



(b) On the basis of the results, you should try 0^{th} order polynomial, 1^{st} order polynomial, 3^{rd} order polynomial and some other order polynomial, show the results include fitting and over-fitting.

1 Transforms Polynomial Feature

```
In [4]: import itertools
import functools
class PolynomialFeature(object):
    """
        polynomial features

        transforms input array with polynomial features

    Example
    ======
    x =
    [[a, b],
     [c, d]]

    y = PolynomialFeatures(degree=2).transform(x)
    y =
    [[1, a, b, a^2, a * b, b^2],
     [1, c, d, c^2, c * d, d^2]]
    """

    def __init__(self, degree=2):
        """
            construct polynomial features

        Parameters
        -----
        degree : int
            degree of polynomial
        """
        assert isinstance(degree, int)
        self.degree = degree

    def transform(self, x):
        """
            transforms input array with polynomial features

        Parameters
        -----
        x : (sample_size, n) ndarray
            input array

        Returns
        -----
        output : (sample_size, 1 + nC1 + ... + nCd) ndarray
            polynomial features
        """
        if x.ndim == 1:
            x = x[:, None]
        x_t = x.transpose()
        features = [np.ones(len(x))]
        for degree in range(1, self.degree + 1):
            for items in itertools.combinations_with_replacement(x_t, degree):
                features.append(functools.reduce(lambda x, y: x * y, items))
        return np.asarray(features).transpose()

class Regression(object):
    """
        Base class for regressors
    """
    pass
```

```
In [5]: # my code

# train data in linear model
def train_linear(x_train,y_train,x_test, degree):
    # feature transform
    ployfeature = PolynomialFeature(degree)
    feature_train = ployfeature.transform(x_train)
    feature_test = ployfeature.transform(x_test)

    # LinearRegression and fit
    linModel = LinearRegression()
    linModel.fit(feature_train,y_train)
    y_pred_test = linModel.predict(feature_test)
    y_pred_train = linModel.predict(feature_train)

    return y_pred_train,y_pred_test
```

2 Regression

2.1 Linear Regression

```
In [6]: class LinearRegression(Regression):
    """
    Linear regression model
    y = X @ w
    t ~ N(t|X @ w, var)
    """

    def fit(self, X:np.ndarray, t:np.ndarray):
        """
        perform least squares fitting

        Parameters
        -----
        X : (N, D) np.ndarray
            training independent variable
        t : (N,) np.ndarray
            training dependent variable
        """
        self.w = np.linalg.pinv(X) @ t
        self.var = np.mean(np.square(X @ self.w - t))

    def predict(self, X:np.ndarray, return_std:bool=False):
        """
        make prediction given input

        Parameters
        -----
        X : (N, D) np.ndarray
            samples to predict their output
        return_std : bool, optional
            returns standard deviation of each predition if True

        Returns
        -----
        y : (N,) np.ndarray
            prediction of each sample
        y_std : (N,) np.ndarray
            standard deviation of each predition
        """
        y = X @ self.w
        if return_std:
```

```

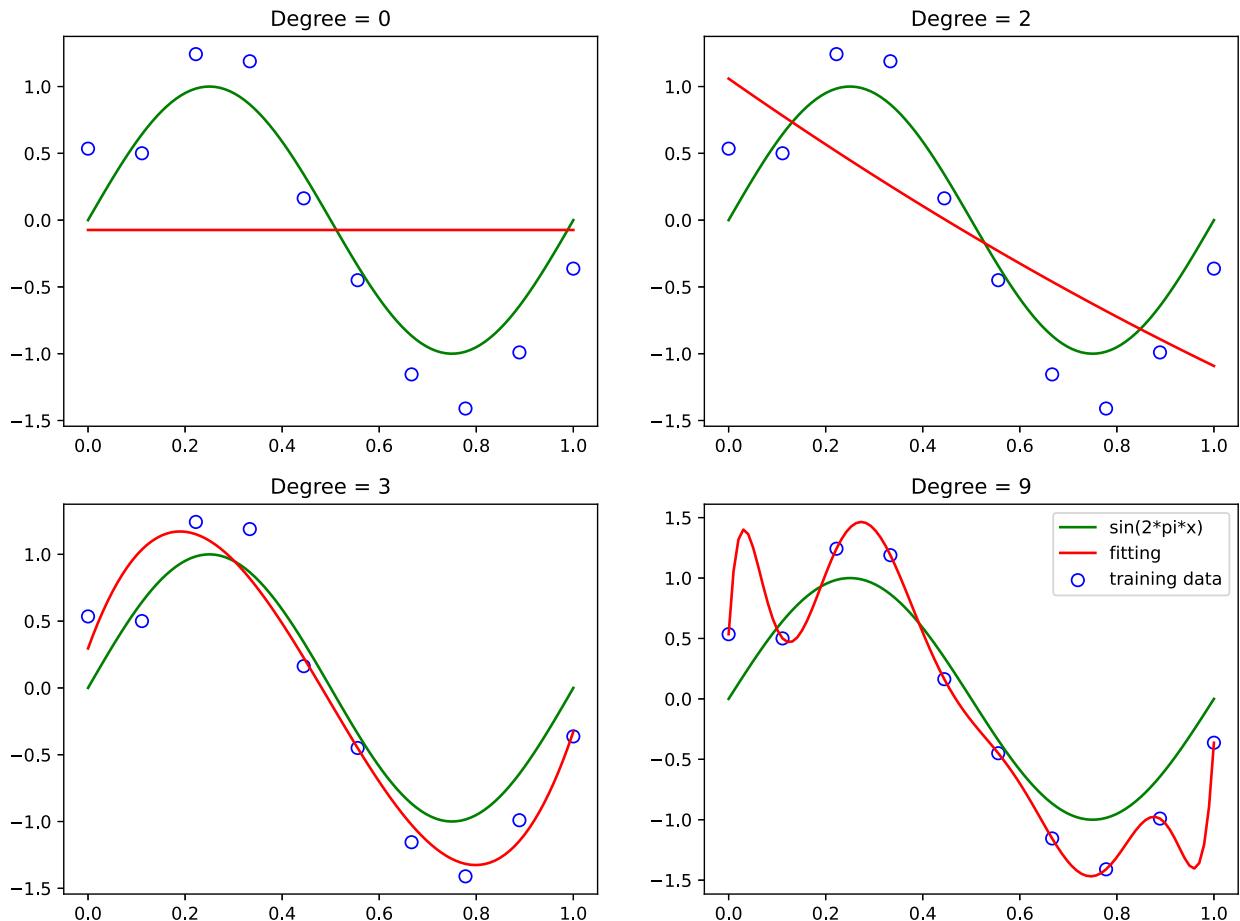
        y_std = np.sqrt(self.var) + np.zeros_like(y)
        return y, y_std
    return y

```

```

In [7]: # train and plot
degree = [0,2,3,9]
idx = 0
plt.figure(figsize=(12,9))
for i in degree:
    y_pred_train,y_pred_test = train_linear(x_train,y_train,x_test, i)
    idx += 1
    plt.subplot(2, 2, idx)
    title = "Degree = " + str(i)
    plt.title(title)
    plt.scatter(x_train, y_train, facecolor="none", edgecolor="b", s=50, label="training data")
    plt.plot(x_test, y_test,c = 'g',label="sin(2*pi*x)")
    plt.plot(x_test, y_pred_test,color='r',label='fitting')
plt.legend()
plt.show()

```



Analysis

We notice that the constant ($M = 0$) and first order ($M = 1$) polynomials give rather poor fits to the data and consequently rather poor representations of the function $\sin(2\pi x)$. The third order ($M = 3$) polynomial seems to give the best fit to the function $\sin(2\pi x)$ of the examples shown in Figure 1.4. When we go to a much higher order polynomial ($M = 9$), we obtain an excellent fit to the training data. In fact, the polynomial passes exactly through each data point and $E(w^*)=0$. However, the fitted curve oscillates wildly and gives a very poor representation of the function $\sin(2\pi x)$. This latter behaviour is known as over-fitting.

2.2 root-mean-square error

(c) Plot the graph of the root-mean-square error.

In [8]:

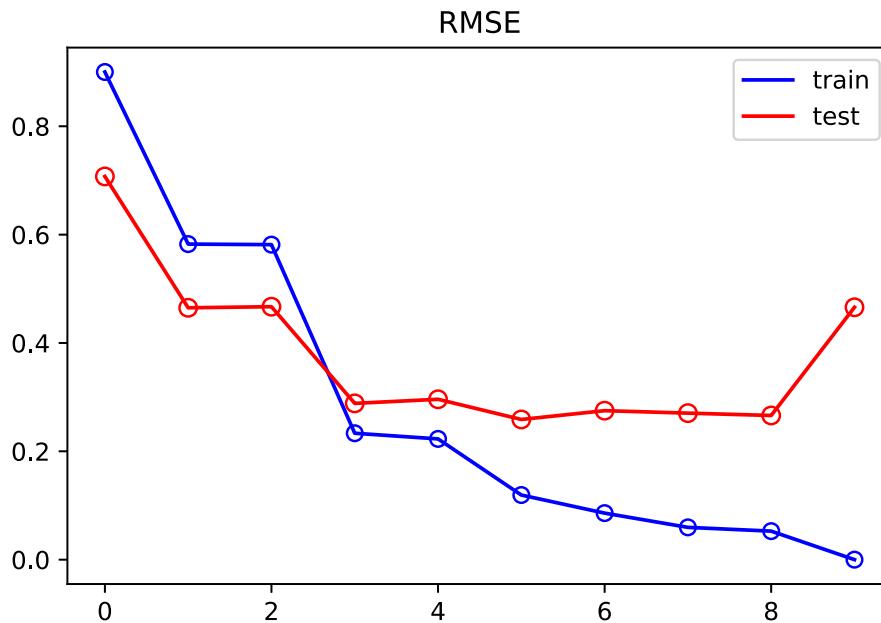
```
# Complete this function
from math import sqrt
from sklearn.metrics import mean_squared_error
def rmse(a, b):
    MSE=mean_squared_error(a,b)
    RMSE=sqrt(MSE)
    return RMSE
```

In [9]:

```
# RMSE
RMSE_train = []
RMSE_test = []
degree = range(0,10)
for i in degree:
    y_pred_train,y_pred_test = train_linear(x_train,y_train,x_test, i)
    RMSE_train.append(rmse(y_pred_train,y_train))
    RMSE_test.append(rmse(y_pred_test,y_test))

plt.title("RMSE")
plt.plot(degree,RMSE_train,color='b',label='train')
plt.plot(degree,RMSE_test,color='r',label='test')
plt.legend()

plt.scatter(degree, RMSE_train, s=40, edgecolors="b", c='', marker='o')
plt.scatter(degree, RMSE_test, s=50, edgecolors="r", c='', marker='o')
plt.show()
```



Analysis

RMS function in which the division by N allows us to compare different sizes of data sets on an equal footing, and the square root ensures that ERMS is measured on the same scale (and in the same units) as the target variable t. Graphs of the training and test set RMS errors are shown, for various values of M, in Figure 1.5. The test set error is a measure of how well we are doing in predicting the values of t for new data observations of x. We note from Figure 1.5 that small values of M give relatively large values of the test set error, and this can be attributed to the fact that the corresponding polynomials are rather inflexible and are incapable of capturing the oscillations in the function $\sin(2\pi x)$.

2.3 Bayesian Regression

(d) Plot the graph of the predictive distribution resulting from a Bayesian treatment of polynomial curve fitting using an M=9 polynomial, with the fixed parameters $\alpha = 5 \times 10^{-3}$ and $\beta = 11.1$ (corresponding to the known noise variance).

```
In [10]: class BayesianRegression(Regression):
    """
        Bayesian regression model

    w ~ N(w|0, alpha^(-1)I)
    y = X @ w
    t ~ N(t|X @ w, beta^(-1))
    """

    def __init__(self, alpha:float=1., beta:float=1.):
        self.alpha = alpha
        self.beta = beta
        self.w_mean = None
        self.w_precision = None

    def _is_prior_defined(self) -> bool:
        return self.w_mean is not None and self.w_precision is not None

    def _get_prior(self, ndim:int) -> tuple:
        if self._is_prior_defined():
            return self.w_mean, self.w_precision
        else:
            return np.zeros(ndim), self.alpha * np.eye(ndim)

    def fit(self, X:np.ndarray, t:np.ndarray):
        """
            bayesian update of parameters given training dataset

        Parameters
        -----
        X : (N, n_features) np.ndarray
            training data independent variable
        t : (N,) np.ndarray
            training data dependent variable
        """

        mean_prev, precision_prev = self._get_prior(np.size(X, 1))

        w_precision = precision_prev + self.beta * X.T @ X
        w_mean = np.linalg.solve(
            w_precision,
            precision_prev @ mean_prev + self.beta * X.T @ t
        )
        self.w_mean = w_mean
        self.w_precision = w_precision
        self.w_cov = np.linalg.inv(self.w_precision)

    def predict(self, X:np.ndarray, return_std:bool=False, sample_size:int=None):
        """
            return mean (and standard deviation) of predictive distribution

        Parameters
        -----
        X : (N, n_features) np.ndarray
            independent variable
        return_std : bool, optional
            flag to return standard deviation (the default is False)
```

```

sample_size : int, optional
    number of samples to draw from the predictive distribution
    (the default is None, no sampling from the distribution)

Returns
-----
y : (N,) np.ndarray
    mean of the predictive distribution
y_std : (N,) np.ndarray
    standard deviation of the predictive distribution
y_sample : (N, sample_size) np.ndarray
    samples from the predictive distribution
"""

if sample_size is not None:
    w_sample = np.random.multivariate_normal(
        self.w_mean, self.w_cov, size=sample_size
    )
    y_sample = X @ w_sample.T
    return y_sample
y = X @ self.w_mean
if return_std:
    y_var = 1 / self.beta + np.sum(X @ self.w_cov * X, axis=1)
    y_std = np.sqrt(y_var)
    return y, y_std
return y

```

```

In [11]: # Write your codes here.
## train data in linear model
def train_bayesian(x_train,y_train,x_test, degree):
    # feature transform
    ployfeature=PolynomialFeature(degree)
    feature_train=ployfeature.transform(x_train)
    feature_test=ployfeature.transform(x_test)

    # BayesianRegression and fit
    BayesianModel = BayesianRegression(alpha=5e-3, beta=11.1)
    BayesianModel.fit(feature_train, y_train)
    y_pred_test= BayesianModel.predict(feature_test)
    y_pred_train = BayesianModel.predict(feature_train)

    # var_train = linModel.var
    # var_test = np.mean(np.square(y_pred - y_test))

    return y_pred_train,y_pred_test

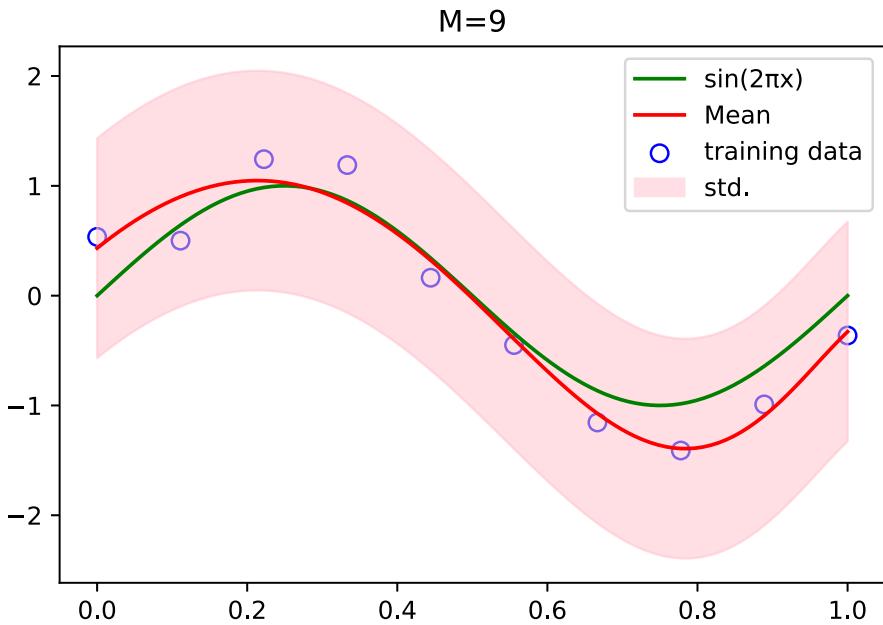
```

```

In [12]: y_pred_train,y_pred_test = train_bayesian(x_train,y_train,x_test, 9)

plt.scatter(x_train, y_train, facecolor="none", edgecolor="b", s=50, label="training data")
plt.plot(x_test, y_test, c="g", label="sin(2πx)")
plt.plot(x_test, y_pred_test, c="r", label="Mean")
plt.fill_between(x_test, y_pred_test-1, y_pred_test+1, color="pink", label="std.", alpha=0.5)
plt.title("M=9")
plt.legend(loc='best')
plt.show()

```



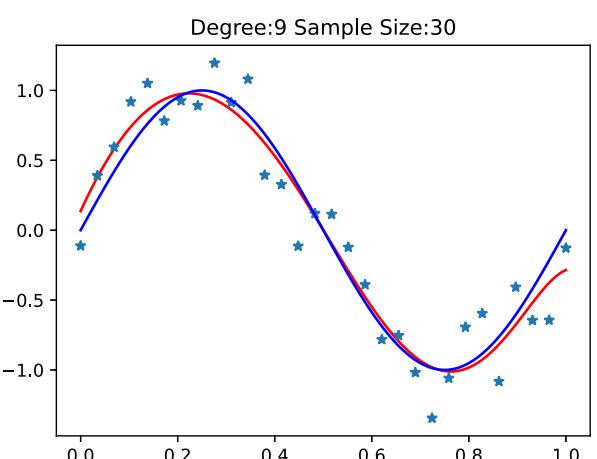
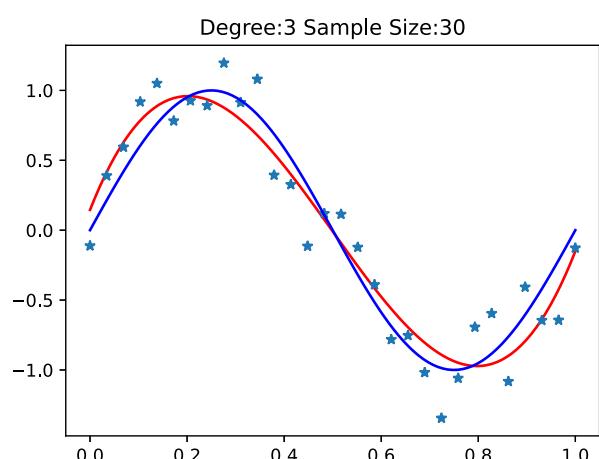
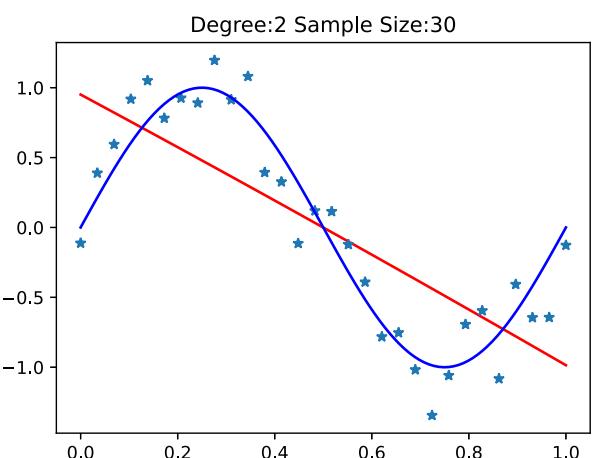
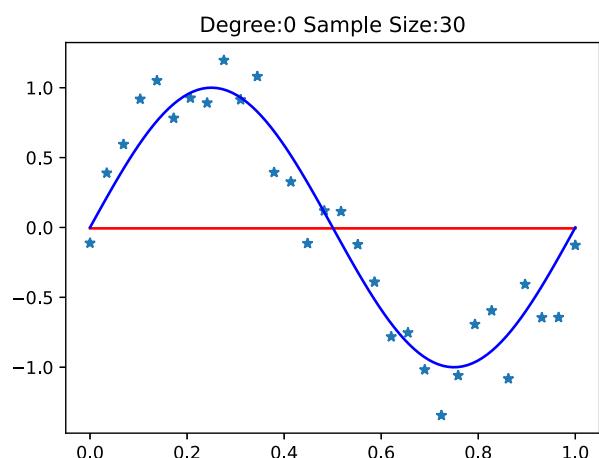
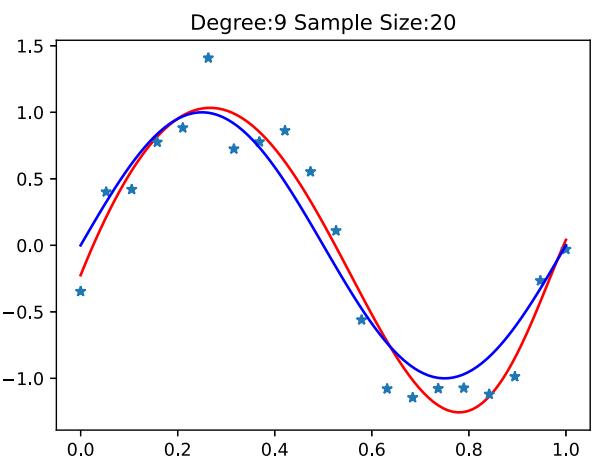
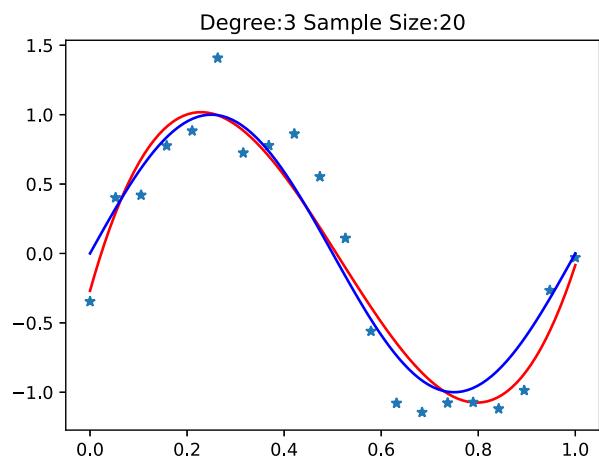
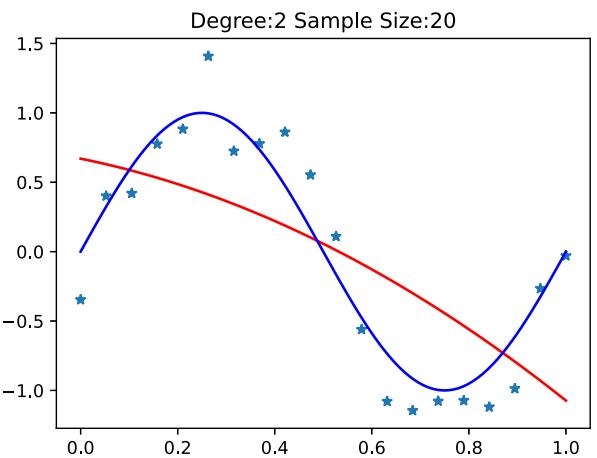
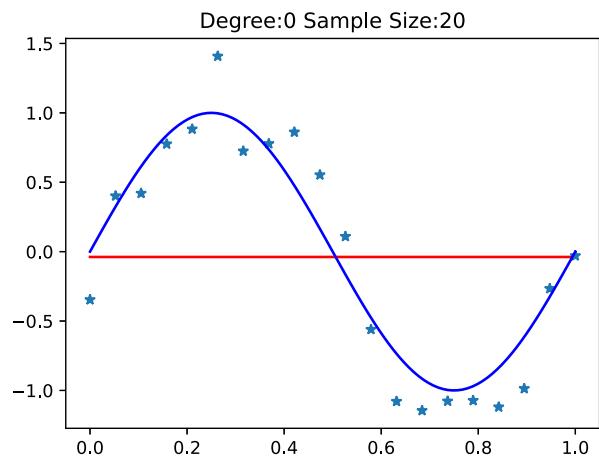
Analysis

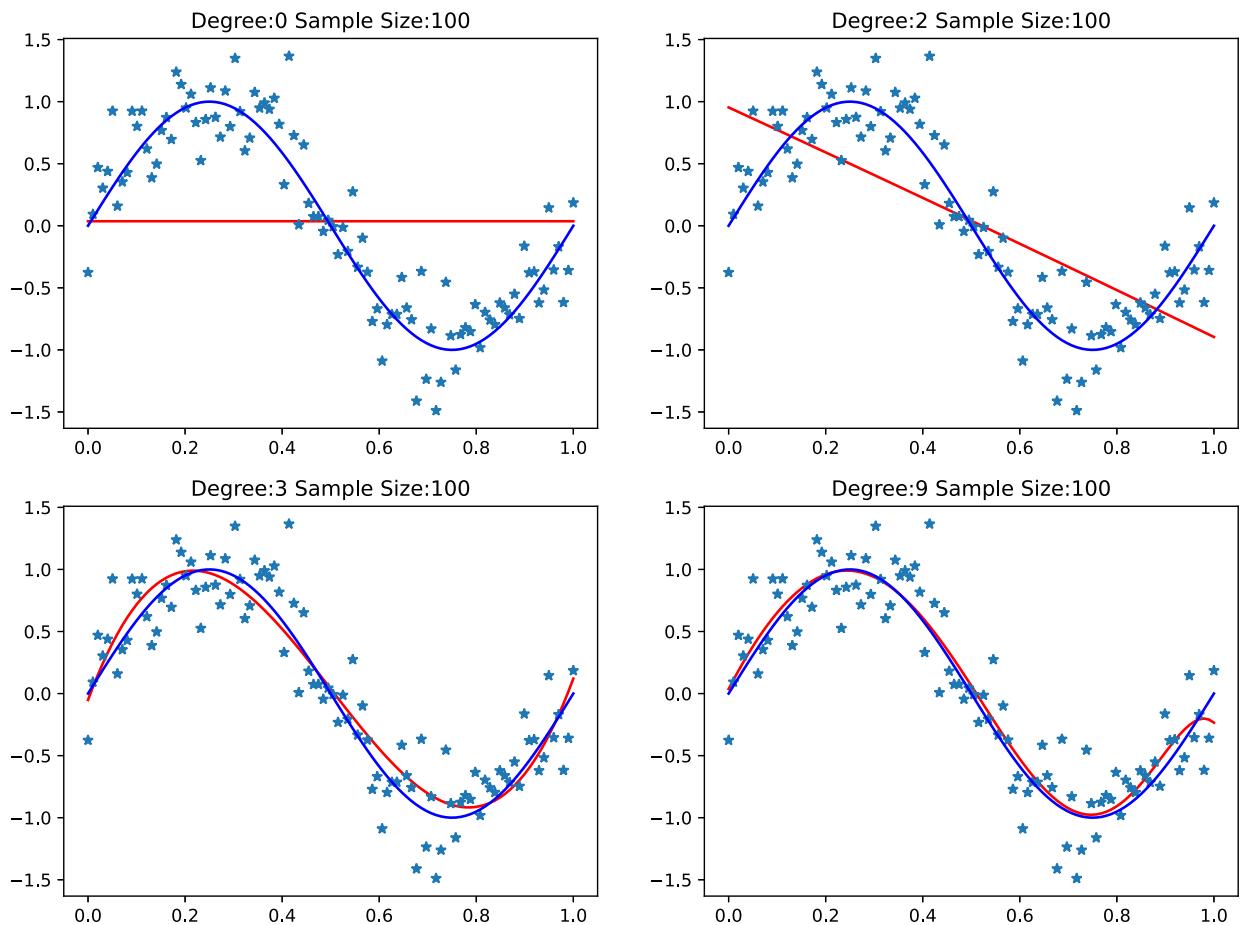
By adopting a Bayesian approach, the over-fitting problem can be avoided. We shall see that there is no difficulty from a Bayesian perspective in employing models for which the number of parameters greatly exceeds the number of data points. Indeed, in a Bayesian model the effective number of parameters adapts automatically to the size of the data set.

2.4 Change the *sample_size*

(e) Change the *sample_size* to 2, 3 or 10 times than before, explain the change of M .

```
In [13]: # Write your codes here.
for j in [2,3,10]:
    x_train, y_train = create_toy_data(func, 10*j, 0.25)
    x_test = np.linspace(0, 1, 100)
    y_test = func(x_test)
    idx = 0
    plt.figure(figsize=(12,9))
    for i in [0,2,3,9]:
        idx += 1
        y_pred_train,y_pred_test = train_bayesian(x_train,y_train,x_test, i)
        plt.subplot(2, 2, idx)
        title = "Degree:"+str(i)+" Sample Size:"+ str(10*j)
        plt.title(title)
        plt.plot(x_test, y_pred_test, c="r")
        plt.plot(x_train, y_train,'*')
        plt.plot(x_test, y_test,c='b')
    plt.show()
```





Analysis

One of the most frequent problems in statistical analysis is the determination of the appropriate sample size. One may ask why sample size is so important. The answer to this is that an appropriate sample size is required for validity. If the sample size it too small, it will not yield valid results. An appropriate sample size can produce accuracy of results. Moreover, the results from the small sample size will be questionable.

A sample size that is too large will result in wasting money and time. It is also unethical to choose too large a sample size. There is no certain rule of thumb to determine the sample size. Some researchers do, however, support a rule of thumb when using the sample size.

For example, in regression analysis, many researchers say that there should be at least 10 observations per variable. If we are using three independent variables, then a clear rule would be to have a minimum sample size of 30. Some researchers follow a statistical formula to calculate the sample size.

Reference

- [1] Kelley, Ken & Maxwell, Scott. (2003). Sample Size for Multiple Regression: Obtaining Regression Coefficients That Are Accurate, Not Simply Significant. *Psychological methods*. 8. 305-21. 10.1037/1082-989X.8.3.305.
- [2] Maxwell, Scott. (2001). Sample size and multiple regression. *Psychological methods*. 5. 434-58. 10.1037//1082-989X.5.4.434.