

## Part 1: The James-Stein Estimator (20 P)

Let  $x_1, \dots, x_N \in \mathbb{R}^d$  be independent draws from a multivariate Gaussian distribution with mean vector  $\mu$  and covariance matrix  $\Sigma = \sigma^2 I$ . It can be shown that the maximum-likelihood estimator of the mean parameter  $\mu$  is the empirical mean given by:

$$\hat{\mu}_{\text{ML}} = \frac{1}{N} \sum_{i=1}^N x_i$$

Maximum-likelihood appears to be a strong estimator. However, it was demonstrated that the following estimator

$$\hat{\mu}_{JS} = \left(1 - \frac{(d-2) \cdot \frac{\sigma^2}{N}}{\|\hat{\mu}_{\text{ML}}\|^2}\right) \hat{\mu}_{\text{ML}}$$

(a shrunk version of the maximum-likelihood estimator towards the origin) has actually a smaller distance from the true mean when  $d \geq 3$ . This however assumes knowledge of the variance of the distribution for which the mean is estimated. This estimator is called the James-Stein estimator. While the proof is a bit involved, this fact can be easily demonstrated empirically through simulation. This is the object of this exercise.

The code below draws ten 50-dimensional points from a normal distribution with mean vector  $\mu = (1, \dots, 1)$  and covariance  $\Sigma = I$ .

```
import numpy

def getdata(seed):

    n = 10          # data points
    d = 50          # dimensionality of data
    m = numpy.ones([d]) # true mean
    s = 1.0         # true standard deviation

    rstate = numpy.random.mtrand.RandomState(seed)
    X = rstate.normal(0,1,[n,d])*s+m

    return X,m,s
```

The following function computes the maximum likelihood estimator from a sample of the data assumed to be generated by a Gaussian distribution:

```
def ML(X):
    return X.mean(axis=0)
```

$$\hat{\mu}_{JS} = \left(1 - \frac{(d-2) \cdot \frac{\sigma^2}{N}}{\|\hat{\mu}_{\text{ML}}\|^2}\right) \hat{\mu}_{\text{ML}}$$

## Implementing the James-Stein Estimator (10 P)

- Based on the ML estimator function, write a function that receives as input the data  $(X_i)_{i=1}^n$  and the (known) variance  $\sigma^2$  of the generating distribution, and computes the James-Stein estimator

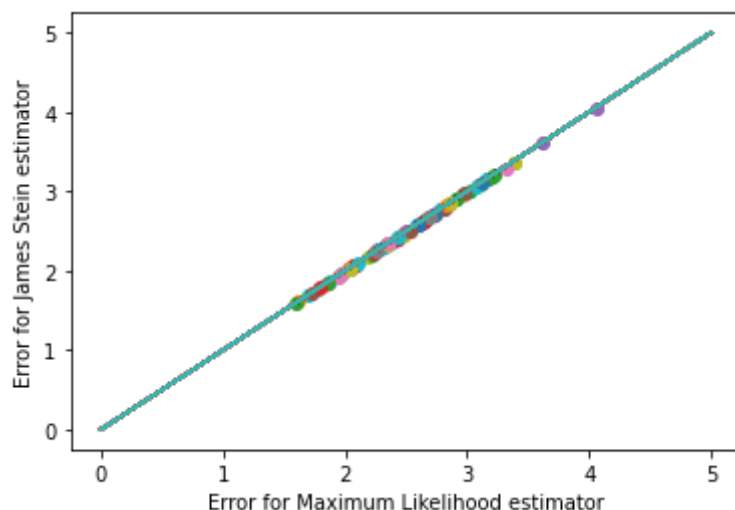
```
def JS(X,s):  
    return ML(X) * (1 - (X.shape[0] - 2) * s * s / X.shape[1] /  
    np.power(np.linalg.norm(ML(X)),2))
```

## Comparing the ML and James-Stein Estimators (10 P)

We would like to compute the error of the maximum likelihood estimator and the James-Stein estimator for 100 different samples (where each sample consists of 10 draws generated by the function `getdata` with a different random seed). Here, for reproducibility, we use seeds from 0 to 99. The error should be measured as the Euclidean distance between the true mean vector and the estimated mean vector.

- Compute the maximum-likelihood and James-Stein estimations.
- Measure the error of these estimations.
- Build a scatter plot comparing these errors for different samples.

```
%matplotlib inline  
import numpy as np  
from matplotlib import pyplot as plt  
  
for seed in range(0, 100):  
    X,m,s = getdata(seed)  
    ML_error = 0.5 * np.sum(np.power(m - ML(X), 2))  
    JS_error = 0.5 * np.sum(np.power(m - JS(X, s), 2))  
    plt.scatter(ML_error, JS_error)  
    plt.plot(np.arange(0, 6), np.arange(0, 6))  
plt.xlabel('Error for Maximum Likelihood estimator')  
plt.ylabel('Error for James Stein estimator')  
plt.show()
```



## Part 2: Bias/Variance Decomposition (30 P)

In this part, we would like to implement a procedure to find the bias and variance of different predictors. We consider one for regression and one for classification. These predictors are available in the module `utils`.

- `utils.ParzenRegressor`: A regression method based on Parzen window. The hyperparameter corresponds to the scale of the Parzen window. A large scale creates a more rigid model. A small scale creates a more flexible one.
- `utils.ParzenClassifier`: A classification method based on Parzen window. The hyperparameter corresponds to the scale of the Parzen window. A large scale creates a more rigid model. A small scale creates a more flexible one. Note that instead of returning a single class for a given data point, it outputs a probability distribution over the set of possible classes.

Each class of predictor implements the following three methods:

- `__init__(self, parameter)`: Create an instance of the predictor with a certain scale parameter.
- `fit(self, X, T)`: Fit the predictor to the data (a set of data points `X` and targets `T`).
- `predict(self, X)`: Compute the output values arbitrary inputs `X`.

To compute the bias and variance estimates, we require *multiple samples* from the training set for a single set of observation data. To accomplish this, we utilize the `Sampler` class provided. The sampler is initialized with the training data and passed to the method for estimating bias and variance, where its function `sampler.sample()` is called repeatedly in order to fit multiple models and create an ensemble of prediction for each test data point.

## Regression Case (15 P)

For the regression case, Bias, Variance and Error are given by:

- $\text{Bias}(Y)^2 = (\mathbb{E}_Y [Y - T])^2$
- $\text{Var}(Y) = \mathbb{E}_Y [(Y - \mathbb{E}_Y[Y])^2]$
- $\text{Error}(Y) = \mathbb{E}_Y [(Y - T)^2]$

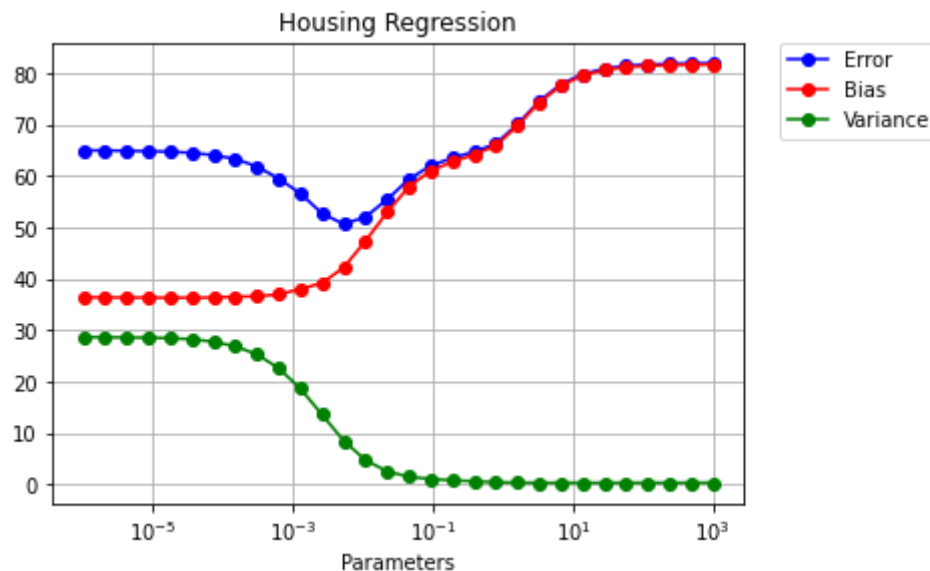
**Task:** Implement the KL-based Bias-Variance Decomposition defined above. The function should repeatedly sample training sets from the sampler (as many times as specified by the argument `nbsamples`), learn the predictor on them, and evaluate the variance on the out-of-sample distribution given by `X` and `T`.

```
def biasVarianceRegression(sampler, predictor, X, T, nbsamples):
    # Train the predictor
    predicted=np.array([predictor.fit(*sampler.sample()).predict(X) for _ in
range(nbsamples)])
    mean = np.mean(predicted, axis=0)

    bias = np.mean(np.power(mean - T, 2))
    variance = np.mean(np.mean(np.power(predicted-mean, 2)))
    return bias, variance
```

Your implementation can be tested with the following code:

```
import utils, numpy
%matplotlib inline
utils.plotBVE(utils.Housing, numpy.logspace(-6, 3, num=30), utils.ParzenRegressor, bi
asVarianceRegression, 'Housing Regression')
```



## Classification Case (15 P)

We consider here the Kullback-Leibler divergence as a measure of classification error, as derived in the exercise, the Bias, Variance decomposition for such error is:

- $\mathrm{Bias}(Y) = D_{\mathrm{KL}}(T \parallel R)$
- $\mathrm{Var}(Y) = \mathbb{E}_Y[D_{\mathrm{KL}}(R \parallel Y)]$
- $\mathrm{Error}(Y) = \mathbb{E}_Y[D_{\mathrm{KL}}(T \parallel Y)]$

where  $R$  is the distribution that minimizes its expected KL divergence from the estimator of probability distribution  $Y$  (see the theoretical exercise for how it is computed exactly), and where  $T$  is the target class distribution.

**Task:** Implement the KL-based Bias-Variance Decomposition defined above. The function should repeatedly sample training sets from the sampler (as many times as specified by the argument `nbsamples`), learn the predictor on them, and evaluate the variance on the out-of-sample distribution given by  $X$  and  $T$ .

```
def KL_div(p, q, axis):
    return np.sum(p * np.log(p / q), axis = axis)
```

```
def biasVarianceClassification(sampler, predictor, X, T, nbsamples=25):

    predicted=np.array([predictor.fit(*sampler.sample()).predict(X) for _ in
range(nbsamples)])

    R = np.exp(np.mean(np.log(predicted), axis=0)) /
np.sum(np.exp(np.mean(np.log(predicted), axis=0)), axis = 1)[:,np.newaxis]

    bias = np.mean(KL_div(T, R, axis=1))
    variance = np.mean(KL_div(R, predicted, axis=2))

    return bias, variance
```

Your implementation can be tested with the following code:

```
import utils,numpy
%matplotlib inline
utils.plotBVE(utils.Yeast,numpy.logspace(-6,3,num=30),utils.ParzenClassifier,bias
sVarianceClassification,'Yeast Classification')
```

