

✓ Notebook 1: Statistical learning. Variance vs Bias

Overview

In this notebook, we will get our hands dirty trying to gain intuition about why machine learning is difficult.

Our task is going to be a simple one, fitting data with polynomials of different order. Formally, this goes under the name of polynomial regression. Here we will do a series of exercises that are intended to give the reader intuition about the major challenges that any machine learning algorithm faces.

Learning Goal

We will explore how our ability to predict depends on the number of data points we have, the "noise" in the data, and our knowledge about relevant features. The goal is to build intuition about why prediction is difficult and discuss general strategies for overcoming these difficulties.

The Prediction Problem

Consider a probabilistic process that gives rise to labeled data (x, y) . The data is generated by drawing samples from the equation

$$y_i = f(x_i) + \eta_i,$$

where $f(x_i)$ is some fixed, but (possibly unknown) function, and η_i is a Gaussian, uncorrelated noise variable such that

$$\langle \eta_i \rangle = 0$$

$$\langle \eta_i \eta_j \rangle = \delta_{ij} \sigma$$

We will refer to the $f(x_i)$ as the **true features** used to generate the data.

To make prediction, we will consider a family of functions $g_\alpha(x; \theta_\alpha)$ that depend on some parameters θ_α . These functions represent the **model class** that we are using to try to model the data and make predictions. The $g_\alpha(x; \theta_\alpha)$ encode the class of **features** we are using to represent the data.

To learn the parameters θ , we will train our models on a **training data set** and then test the effectiveness of the model on a *different* dataset, the **test data set**. The reason we must divide our data into a training and test dataset is that the point of machine learning is to make accurate predictions about new data we have not seen. As we will see below, models that give the best fit to the training data do not necessarily make the best predictions on the test data. This will be a running theme that we will encounter repeatedly in machine learning.

For the remainder of the notebook, we will focus on polynomial regression. Our task is to model the data with polynomials and make predictions about the new data that we have not seen. We will consider two qualitatively distinct situations:

- In the first case, the process that generates the underlying data is in the model class we are using to make predictions. For polynomial regression, this means that the functions $f(x_i)$ are themselves polynomials.
- In the second case, our data lies outside our model class. In the case of polynomial regression, this could correspond to the case where the $f(x_i)$ is a 10-th order polynomial but $g_\alpha(x; \theta_\alpha)$ are polynomials of order 1 or 3.

In the exercises and discussion we consider 3 model classes:

- the case where the $g_\alpha(x; \theta_\alpha)$ are all polynomials up to order 1 (linear models),
- the case where the $g_\alpha(x; \theta_\alpha)$ are all polynomials up to order 3,
- the case where the $g_\alpha(x; \theta_\alpha)$ are all polynomials up to order 10.

To measure our ability to predict, we will learn our parameters by fitting our training dataset and then making predictions on our test data set. One common measure of predictive performance of our algorithm is to compare the predictions, $\{y_j^{\text{pred}}\}$, to the true values $\{y_j\}$. A commonly employed measure for this is the sum of the mean square-error (MSE) on the test set:

$$MSE = \frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} (y_j^{\text{pred}} - y_j)^2$$

We will return to this in later notebooks. For now, we will try to get a qualitative picture by examining plots on test and training data.

Fitting vs. predicting when the data is in the model class

We start by considering the case:

$$f(x) = 2x.$$

Then the data is clearly generated by a model that is contained within all three model classes we are using to make predictions (linear models, third order polynomials, and tenth order polynomials).

Run the code for the following cases:

- For $f(x) = 2x$, $N_{\text{train}} = 10$ and $\sigma = 0$ (noiseless case), train the three classes of models (linear, third-order polynomial, and tenth order polynomial) for a training set when $x_i \in [0, 1]$. Make graphs comparing fits for different order of polynomials. Which model fits the data the best?
- Do you think that the data that has the least error on the training set will also make the best predictions? Why or why not? Can you try to discuss and formalize your intuition? What can go right and what can go wrong?
- Check your answer by seeing how well your fits predict newly generated test data (including on data outside the range you fit on, for example $x \in [0, 1.2]$) using the code below. How well do you do on points in the range of x where you trained the model? How about points outside the original training data set?
- Repeat the exercises above for $f(x) = 2x$, $N_{\text{train}} = 10$, and $\sigma = 1$. What changes?
- Repeat the exercises above for $f(x) = 2x$, $N_{\text{train}} = 100$, and $\sigma = 1$. What changes?
- Summarize what you have learned about the relationship between model complexity (number of parameters), goodness of fit on training data, and the ability to predict well.

Fitting vs. predicting when the data is not in the model class

Thus far, we have considered the case where the data is generated using a model contained in the model class. Now consider $f(x) = 2x - 10x^5 + 15x^{10}$. Notice that for linear and third-order polynomial the true model $f(x)$ is not contained in model class $g_a(x)$.

- Repeat the exercises above fitting and predicting for $f(x) = 2x - 10x^5 + 15x^{10}$ for $N_{\text{train}} = 10, 100$ and $\sigma = 0, 1$. Record your observations.
- Do better fits lead to better predictions?
- What is the relationship between the true model for generating the data and the model class that has the most predictive power? How is this related to the model complexity? How does this depend on the number of data points N_{train} and σ ?
- Summarize what you think you learned about the relationship of knowing the true model class and predictive power.

✓ Training the models:

```
#This is Python Notebook to walk through polynomial regression examples
#We will use this to think about regression
import numpy as np
%matplotlib inline

from sklearn import datasets, linear_model
from sklearn.preprocessing import PolynomialFeatures

from matplotlib import pyplot as plt, rcParams
fig = plt.figure(figsize=(8, 6))

# The Training Data
N_train=1000
sigma_train=0;

# Train on integers
x=np.linspace(0.05,0.95,N_train)
# Draw Gaussian random noise
s = sigma_train*np.random.randn(N_train)

#linear
#y=2*x+s

# Tenth Order
y=2*x-10*x**5+15*x**10+s

p1=plt.plot(x, y, "o", ms=8, alpha=0.5, label='Training')

# Linear Regression : create linear regression object
clf = linear_model.LinearRegression()

# Train the model using the training set
# Note: sklearn requires a design matrix of shape (N_train, N_features). Thus we reshape x to (N_train, 1):
clf.fit(x[:, np.newaxis], y)

# Use fitted linear model to predict the y value:
xplot=np.linspace(0.02,0.98,200) # grid of points, some are in the training set, some are not
linear_plot=plt.plot(xplot, clf.predict(xplot[:, np.newaxis]), label='Linear')
```

```

# Polynomial Regression
poly3 = PolynomialFeatures(degree=3)
# Construct polynomial features
X = poly3.fit_transform(x[:,np.newaxis])
clf3 = linear_model.LinearRegression()
clf3.fit(X,y)

Xplot=poly3.fit_transform(xplot[:,np.newaxis])
poly3_plot=plt.plot(xplot, clf3.predict(Xplot), label='Poly 3')

# Fifth order polynomial in case you want to try it out
#poly5 = PolynomialFeatures(degree=5)
#X = poly5.fit_transform(x[:,np.newaxis])
#clf5 = linear_model.LinearRegression()
#clf5.fit(X,y)

#Xplot=poly5.fit_transform(xplot[:,np.newaxis])
#plt.plot(xplot, clf5.predict(Xplot), 'r--',linewidth=1)

poly10 = PolynomialFeatures(degree=10)
X = poly10.fit_transform(x[:,np.newaxis])
clf10 = linear_model.LinearRegression()
clf10.fit(X,y)

Xplot=poly10.fit_transform(xplot[:,np.newaxis])
poly10_plot=plt.plot(xplot, clf10.predict(Xplot), label='Poly 10')

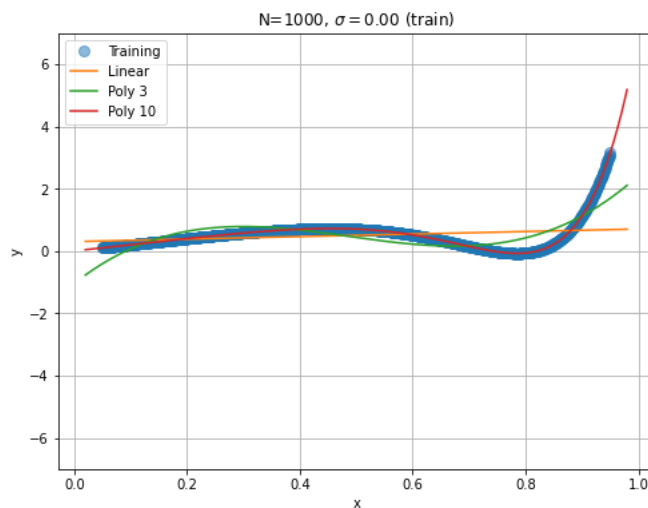
plt.legend(loc='best')
plt.ylim([-7,7])
plt.xlabel("x")
plt.ylabel("y")
Title="N=%i,  $\sigma$ =%.2f"%(N_train,sigma_train)
plt.title(Title+" (train)")

# Linear Filename
filename_train="train-linear_N=%i_noise=%.2f.pdf"%(N_train, sigma_train)

# Tenth Order Filename
#filename_train="train-o10_N=%i_noise=%.2f.pdf"%(N_train, sigma_train)

# Saving figure and showing results
plt.savefig(filename_train)
plt.grid()
plt.show()

```



✓ Testing the fitted models

```

# Generate Test Data
%matplotlib inline
# Number of test data

```

```

N_test=20
sigma_test=sigma_train

# Generate random grid points (x) in the interval [0, max_x]:
# Note some points will be drawn outside the training interval
max_x=1.2
x_test=max_x*np.random.random(N_test)

# Draw random Gaussian noise
s_test = sigma_test*np.random.randn(N_test)

# Linear
y_test=2*x_test+s_test
# Tenth order
#y_test=2*x_test-10*x_test**5+15*x_test**10+s_test

# Make design matrices for prediction
x_plot=np.linspace(0,max_x, 200)
X3 = poly3.fit_transform(x_plot[:,np.newaxis])
X10 = poly10.fit_transform(x_plot[:,np.newaxis])

##### PLOTTING RESULTS #####

fig = plt.figure(figsize=(8, 6))

p1=plt.plot(x_test, y_test, 'o', ms=10, alpha=0.5, label='test data')
p2=plt.plot(x_plot,clf.predict(x_plot[:,np.newaxis]), lw=2, label='linear')
p3=plt.plot(x_plot,clf3.predict(X3), lw=2, label='3rd order')
p10=plt.plot(x_plot,clf10.predict(X10), lw=2, label='10th order')

plt.legend(loc='best')
plt.xlabel('x')
plt.ylabel('y')
plt.legend(loc='best')
Title="N=%i,  $\sigma$ =%.2f"%(N_test,sigma_test)
plt.title(Title+" (pred.)")
plt.tight_layout()
plt.ylim((-6,12))

# Linear Filename
filename_test="pred-linear_N=%i_noise=%.2f.pdf"%(N_test, sigma_test)

# Tenth Order Filename
#filename_test=Title+"pred-o10.pdf"

# Saving figure and showing results
plt.savefig(filename_test)
plt.grid()
plt.show()

```

