

DATA SCIENCE II: Machine Learning MTH 9899 Baruch College

Lecture 1: Introduction to Machine Learning

Adrian Sisser Dmytro Karabash

March 30, 2016

Outline

- 1 General Topics in Machine Learning
 - Supervised vs Unsupervised Learning
 - Classification vs Regression
 - Cross Validation
 - Bias vs Variance
 - Bias vs Variance Example - Ridge Regression
- 2 "Modern" Wave of Machine Learning
 - Logistic Regression
 - Stochastic Gradient Descent and Batch Learning
 - Keras
 - Further Preview of Keras: Multi-Layer Perceptron
 - Further Preview of Keras: Regularizations

Supervised Learning

Supervised machine learning consists of learning a function from a set of labeled training examples.

- Generally, You are given input examples AND output values.
- Success can be easily measured through a variety of metrics on in and out-of-sample observations.
- Sometimes, we don't have exact output values, but instead, a notion of 'maximizing' a function (ie Reinforcement Learning).

Unsupervised Learning

In Unsupervised Learning, you're trying to learn a structure that you don't know at the beginning. There are 2 main categories of Unsupervised Learning:

- Clustering – Identify similar/related items based on their features.
 - Identify 'similar' stocks based on returns or other characteristics.
 - Group mortgages together based on geographic data to understand default correlations.
- Latent Variable Models – Identify underlying variables that drive the features you can observe.
 - Latent Variable - A variable whose value is never known, but instead is implied by its state.

Classification

Identifying which category a variable belongs to. The categories can be:

- **Ordinal** Variables - which have an intrinsic order, ie. credit ratings
- **Categorical** Variables - No implicit ordering, such as what industry a stock belongs to.

Classification Metrics

Receiver Operating Characteristic (ROC) - A graph of the true positive vs false positive rate parameterized by the cutoff used to discriminate between outcomes in a true/false classification.

Confusion Matrix - A table of correct vs incorrect values across all categories.

Regression

Regression refers to prediction a continuous numerical variable - which is what we will focus on in this course. There are a wide variety of different metrics to measure the quality of fit of a regression, each with their own strengths and weaknesses.

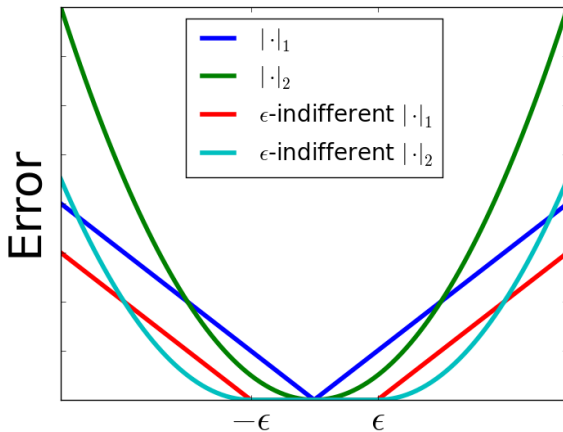
- Mean Squared Error (L2 Norm) - MSE is the most common metric for regression because it's intuitive and very easy to calculate. The problem is, it's not robust to outliers.
- Absolute Error (L1 Norm) - The L1 Norm is a very robust metric that can deal with outliers. Unfortunately, it's very costly to optimize since it's not a convex problem.

Regression

- ϵ Indifferent - This is a metric where we don't penalize for things within some constant, ϵ , of the training value, then apply another metric (ie L1 or L2 norm) to points outside of this area.

Ultimately, the best metric is a tradeoff of computational speed, robustness, and the underlying goal.

Regression



Regression

You might need to do some cleaning or filtering if you're using an L2 Norm metric.

- Median Absolute Deviation (MAD) Filtering - Calculate the MAD - defined as the median of the deviation of every point in a series from the series' median. Then pull in all points to be within n (typically 3 to 5) MADs of the median:

$$\text{MAD}_X = \text{med} |x_i - \text{med } X| \quad (1)$$

$$x'_i = \text{clip}(x_i, \text{med}_X + / - n \text{MAD}_X) \quad (2)$$

- Winsorization - Clip points to a given percentile or number of σ of μ

Summary

- Cross Validation (CV) refers to fitting your model on a portion of your data, and testing it on the out-of-sample portion repeatedly, omitting a different section each time.
- First, divide your data up into F chunks, or 'folds'.
- For example, for the 3-Fold CV shown below, we would fit a model to portions 1,2 and 2,3 and use those models to calculate the error on samples 3 and 4 respectively.
- At an extreme, we can perform 'Leave One Out' CV, which is equivalent to N -Fold CV.
- Model Selection - Does adding a new variable really improve our model?

Motivation

- In ML, we face a risk of overfitting our model to the data. CV will help us avoid this, by measuring the performance on data that is not used to build the underlying model.
- Without CV, we will overestimate the accuracy of our models. This is what leads us to adjusted R^2 .
- Regularization is a valuable technique we will use, and CV is well-suited to calibrating parameters.

Caveats

- With time series data, you have to be careful. If you divide your data up such that a single time is distributed across multiple folds, you might be using forward looking data! For example, imagine you're fitting a complex Neural Network
- Without CV, we will overestimate the accuracy of our models. This is what leads us to adjusted R^2 .
- Regularization is a valuable technique we will use, and CV is well-suited to calibrating parameters.
- Model Selection - Does adding a new variable really improve our model?

Bias & Variance

Normally, we try to calculate unbiased estimators:

$$E[\hat{x}] = x \quad (3)$$

Sometimes, we'd rather introduce a bias, if it can reduce the variance of our predictions. We have an inherent *variance* in our predictor, based on the input dataset, which can be considered a random sample.

Bias vs Variance in Linear Regression

In a traditional linear model (LM), of the form $Y = X\beta + \epsilon$, we assume $\epsilon \sim N(0, \sigma^2)$ and $\hat{\beta} = (X^T X)^{-1} X^T Y$. Let's consider our expected squared prediction error:

$$\begin{aligned}\mathbb{E}_{\hat{\beta}}[\|y - \hat{\beta}x\|_2] &= \mathbb{E}_{\hat{\beta}}[\|(\beta x + \epsilon) - \hat{\beta}x\|_2] \\&= \mathbb{E}_{\hat{\beta}}[\|(\beta x + \epsilon)\|_2 - 2(\beta x + \epsilon)\hat{\beta}x + \|(\beta x + \epsilon)\|_2] \\&= \mathbb{E}_{\hat{\beta}}(\hat{\beta}x)^2 + \epsilon^2 - 2\beta x \mathbb{E}_{\hat{\beta}} \hat{\beta}x + (\beta x)^2 \\&= [\mathbb{E}_{\hat{\beta}}[(\hat{\beta}x)^2] - \mathbb{E}_{\hat{\beta}}[\hat{\beta}x]^2] + \epsilon^2 + (\mathbb{E}_{\hat{\beta}}[\hat{\beta}x] - \beta x)^2\end{aligned}$$

Variance of
the Predictor

Irreducible
error

Bias of the
Predictor

Bias vs Variance in Linear Regression

For normal linear regression, $\mathbb{E} \hat{\beta} = \beta$, so the bias term can be eliminated:

$$\begin{aligned}\mathbb{E}_{\hat{\beta}}[\|y - \hat{\beta}x\|_2] &= [\mathbb{E}_{\hat{\beta}}[(\hat{\beta}x)^2] - \mathbb{E}_{\hat{\beta}}[\hat{\beta}x]^2] + \epsilon^2 + (\mathbb{E}_{\hat{\beta}}[\hat{\beta}x] - \beta x)^2 \\ &= [\mathbb{E}_{\hat{\beta}}[(\hat{\beta}x)^2] - \mathbb{E}_{\hat{\beta}}[\hat{\beta}x]^2] + \epsilon^2\end{aligned}$$

Later, when we cover Ridge Regression, we'll see if allowing a biased estimator can improve the overall prediction quality.

Ridge Regression

Ridge Regression, also known as Tikhonov Regularization, is an extension of a normal least squares regression.

- We will add a penalty term to our normal linear regression - $\lambda \|\beta\|_2$
- This forces a tradeoff between the magnitude of the β s and the error terms.
- This is an example of **Regularization**, the notion of adding a penalty to shrink fitted parameters and reduce variance.

Ridge Regression

$$\min_{\beta^R} \|Y - X\hat{\beta}^R\| + \lambda\|\beta^R\|_2$$

$$\mathcal{L} = \|Y - X\hat{\beta}^R\| + \lambda\|\hat{\beta}^R\|_2$$

$$\frac{\partial \mathcal{L}}{\partial \hat{\beta}^R} = -2Y^T X + 2X^T X \hat{\beta}^R + 2\lambda \hat{\beta}^R$$

$$\hat{\beta}^R = (X^T X + \lambda I)^{-1} X^T Y$$

Regularization
penalty

Ridge Regression

As we saw, in Ridge Regression, the estimate of β is no longer unbiased, in fact, we can show that:

$$\begin{aligned}\text{Bias}(\hat{\beta}^R) &= -\lambda(X^T X + \lambda I)^{-1} \beta \\ \text{Var}(\hat{\beta}^R) &= \sigma^2(X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}\end{aligned}$$

Compare this to OLS:

$$\begin{aligned}\text{Bias}(\hat{\beta}^{LS}) &= 0 \\ \text{Var}(\hat{\beta}^{LS}) &= \sigma^2(X^T X)^{-1}\end{aligned}$$

Ridge Regression

Does it work? Does allowing a biased estimator with a lower variance improve our regression results? Let's see in an IPython Notebook.

Ridge Regression - Choosing λ

So how do we pick a good λ ?? Cross-Validation!!

- We test various values of λ and use the value that minimizes the out-of-sample MSE.
- For Ridge, we can do something called 'Generalized Cross Validation', which is an estimate of leave-one-out validation

Ridge Regression

A few caveats:

- Ridge Regression is dependent on scale. Since we penalize all β values equally, we need to make sure that all variables are normalized, ie $\mu = 0, \sigma = 1$.
- The λ value scales linearly with the number of points, ie if your sample size doubles, your λ should too.

$$\hat{\beta}^R = (X^T X + \lambda I)^{-1} X^T Y$$

- Logistic regression is the most basic model for conditional distribution of some variable y given some x .
- It models probability via linear model that is transformed via sigmoid

$$\frac{1}{1 + e^{\beta_0 + \beta_1 x_1}}$$

- If one has to predict either 0 or 1 rather than probability $\beta_0 + \beta_1 x_1$ is linear classifier.
- Simplicity is perhaps the main reasons why Logistic Regression is so widely used in discrete data analysis (see Boltzmann distribution with two states for more details).

- Stochastic gradient descent is going down the gradient example by example since typically the function we are trying to minimize has form $\sum_i f(e_i)$ where sum is over examples.
- Hence the change is done via subtracting $\eta \nabla f(e_i)$, where η is the learning rate.
- Batch learning has the same idea except examples are grouped in *batches*, for example of size 100, and then for butch K we subtract $\eta \sum_{i=K+1}^{100*K+100} \nabla f(e_i)$.
- *This does not make sense for all models specifically for models with closed solutions, but for non-linear models this technique had a lot of success especially if f has some properties like convexity.*


```
model = Sequential()  
model.add(Dense(20, init='zero', Activation='sigmoid'))  
sgd = SGD(lr=0.1, decay=1e-6,  
          momentum=0.9, nesterov=True)  
model.compile(loss='binary_crossentropy', optimizer=sgd)
```

```
model = Sequential()  
model.add(Dense(64, input_dim=20, activation='relu'))  
model.add(Dropout(0.5))  
model.add(Dense(64, activation='relu'))  
model.add(Dropout(0.5))  
model.add(Dense(10, activation='softmax'))  
model.compile(  
    loss='categorical_crossentropy', optimizer='adadelata')
```

We wrap this together by introducing simple 4-line code for logistic regression.

```
model = Sequential()  
model.add(Dense(64, input_dim=20,  
    W_regularizer=l2(0.01)), activation='relu'))  
model.add(Dropout(0.5))  
model.add(Dense(64,  
    W_regularizer=l1(l1=0.01), activation='relu'))  
model.add(Dropout(0.5))  
model.add(Dense(10,  
    W_regularizer=l1l2(l1=0.01, l2=0.01),  
    activation='softmax'))  
model.compile(  
    loss='categorical_crossentropy', optimizer='adadelta')
```