# Learner Implementation

Tuesday, June 7, 2016   2:00 PM

Constructor parameters:

1. Lin Reg Learner ()
2. Polynomial Learner (degree)
3. KNN Learner (K)
4. Kernel Regression Learner (K, weighting func)
   (basically a Weighed KNN)
5. Decision Tree Learner (threshold value, max_depth, min_split/max_leaf_size

1. **LinRegLearner.py**

```
class LinRegLearner::
    def __init__():
        pass
    def addEvidence(Xtrain, Ytrain):
        self.m, self.b = (yflrf).linreg(Xtrain, Ytrain)
    def query(Xtest):
        return (self.m * Xtest + self.b)
```

**Linear**
We already know the result should be in form of
y = mx + b

Loop until delta_J < threshold:
    sse = 0
    for x, y in zip(XTrain, YTrain):
        ypred = mx + b
        sse += (ypred - y)**2

    J = 1/2 * sse

    m_new = m - (ypred-y)*x*LearningRate
    b_new = b - (ypred - y)*LearningRate

API:

learner = learner()
Learner.addEvidence(Xtrain, Ytrain)
y = learner.query(Xtest)

J is the function I want to minimize here, as the 'loss function'

If we do the derivative, which is:
dJ/dm = (ypred - y)x
dJ/db = (ypred - y)

Without normalization, oscillation around optimum point is possible, therefore we need to adjust by learning rate.

We can set a threshold that if change of J is smaller than it we can quit the program.
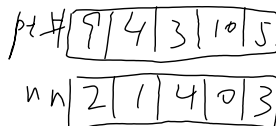
2. **KNNLearner.py**

```
Class KNNLearner::
    def __init__(K=3):
        self.K = K
    def addEvidence(Xtrain, Ytrain):
        self.Xtrain = Xtrain
        self.Ytrain = Ytrain
    def query(Xtest):
        for i in range(0, self.Xtrain.shape[0]):
            dist[i] = metric(self.Xtrain[i], Xtest)
        nn = dist.argsort()

    for idx in nn:
        ysum += Ytrain[idx]

    return ysum/K
```
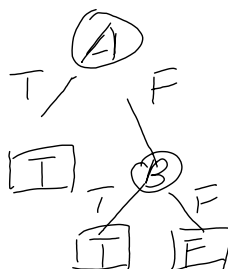


In a classification problem, the only difference lays here. We need a Votes to give classification. E.g.
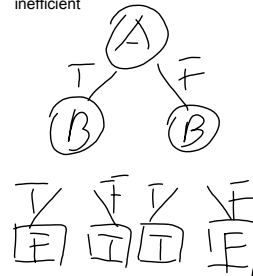
Votes[Ytrain[idx]] += 1

return (the most voted class)

**About Project 1**

Pandas indexing change

Before 0.15, df[0:2] works as numpy index inside
    df[0] gives the first row
After 0.15, df[0] works differently returns that index equals to 0. df.ix[0] gives the same results as older version





3. **Decision Tree: CART**

Bad: hard to learn "optimal decision tree"
1. Greedy
2. Depth/Size
3. NP-complete (cannot be proven to be computable in polynomial time)
4. Degenerate cases

Good:
1. Easy to understand
2. No preprocessing of data
3. Do both classification and regression
4. 'White box' model

'A or B' example



'A Xor B' case would be a problem, we have to create complete tree; if factors increases, we need 2^N nodes; becomes inefficient