

MIE 1624

In class presentation: Linear Regression

Group 7:

Akshay Vishwakarma

Haohan Li

Parth Dave

Rui Fang

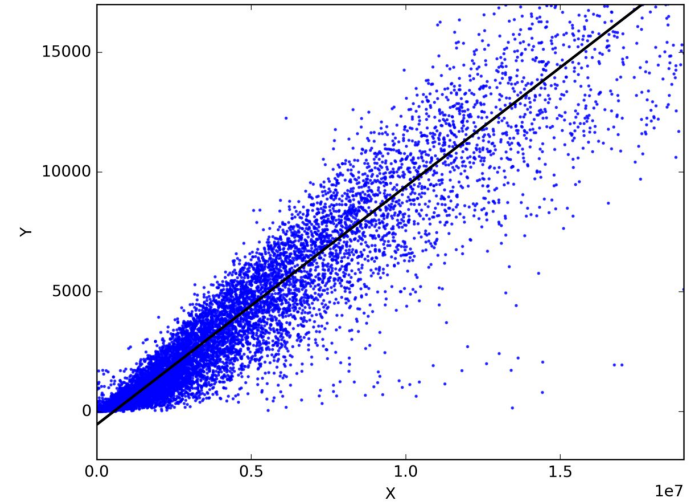
Xi Sun





What is Linear Regression ?

- Prediction tool that helps in predicting an outcome (dependent) variable
- Also helps identify which variable are significant predictors of the outcome variable
 - Indicated by the magnitude and sign of the variable

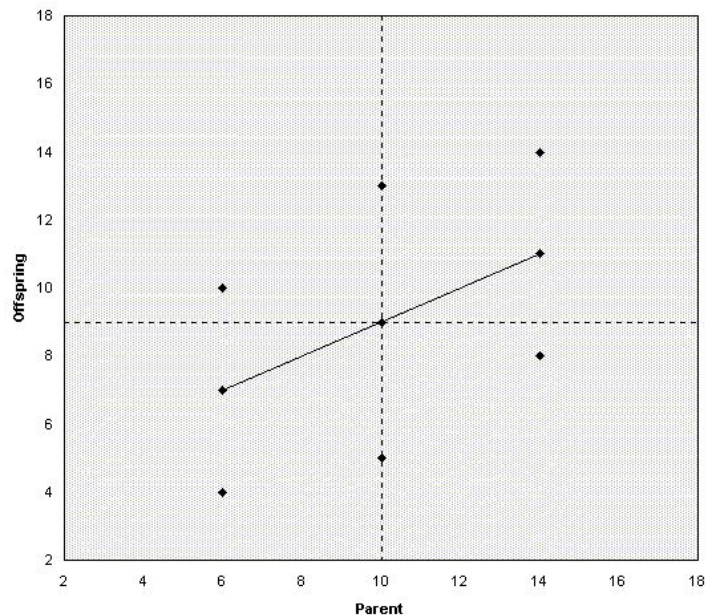




History behind Linear Regression

Originally invented by Sir Francis Galton in 1875 for the Genetic study.

- First model was pea size study parent pea plant and the offspring
- Later the hereditary problem study used the same model





Objective of Linear Regression

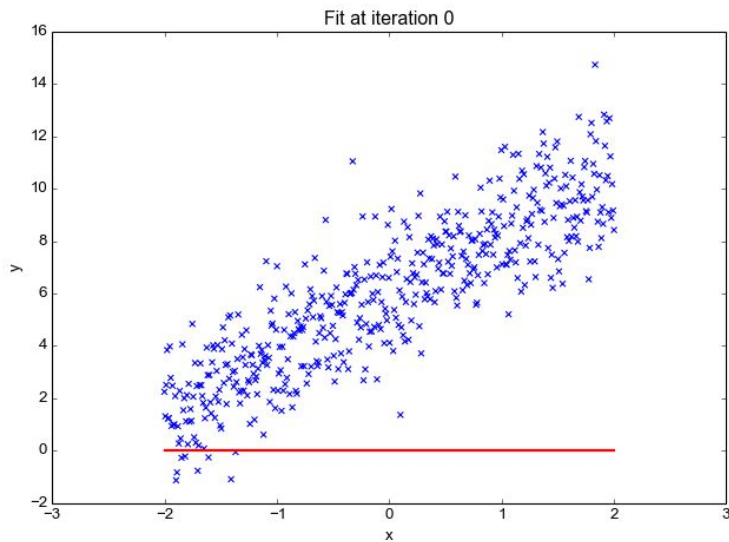
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

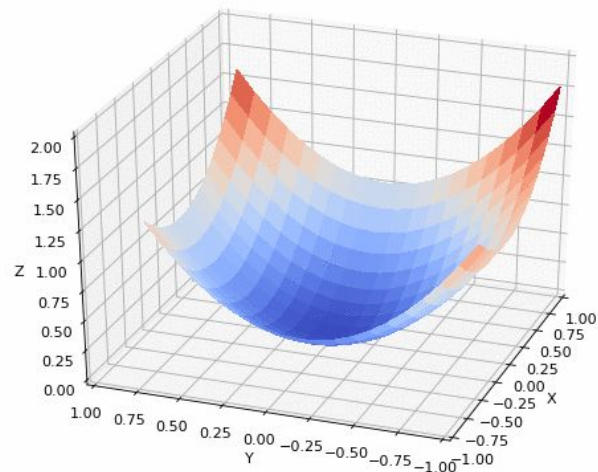
Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Linear Regression in action

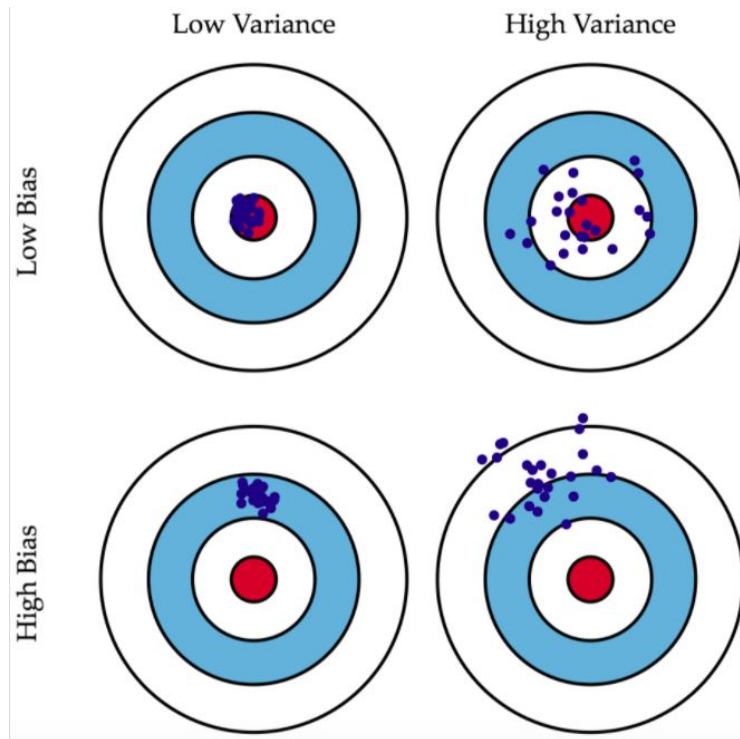


Fitting the linear model



Cost function minimization

Bias-Variance tradeoff





Bias-Variance tradeoff

$$y = f(x) + \varepsilon$$

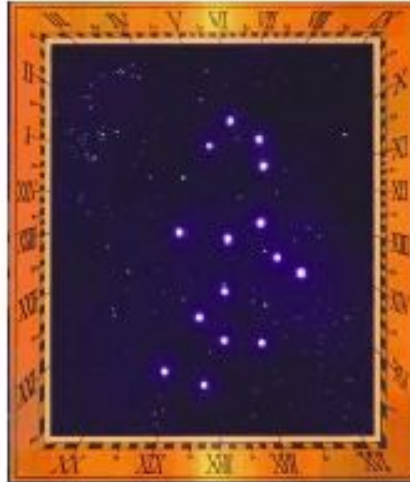
$$\begin{aligned} \mathbb{E}[(y - \hat{f})^2] &= \mathbb{E}[y^2 + \hat{f}^2 - 2y\hat{f}] \\ &= \mathbb{E}[y^2] + \mathbb{E}[\hat{f}^2] - \mathbb{E}[2y\hat{f}] \\ &= \text{Var}[y] + \mathbb{E}[y]^2 + \text{Var}[\hat{f}] + \mathbb{E}[\hat{f}]^2 - 2f\mathbb{E}[\hat{f}] \\ &= \text{Var}[y] + \text{Var}[\hat{f}] + (f^2 - 2f\mathbb{E}[\hat{f}] + \mathbb{E}[\hat{f}]^2) \\ &= \text{Var}[y] + \text{Var}[\hat{f}] + (f - \mathbb{E}[\hat{f}])^2 \\ &= \sigma^2 + \text{Var}[\hat{f}] + \text{Bias}[\hat{f}]^2 \end{aligned}$$

Overfitting

data



data



data



Overfitting

normal



normal



normal





Overfitting

overfitiing



overfitiing



overfitiing





Overfitting

data

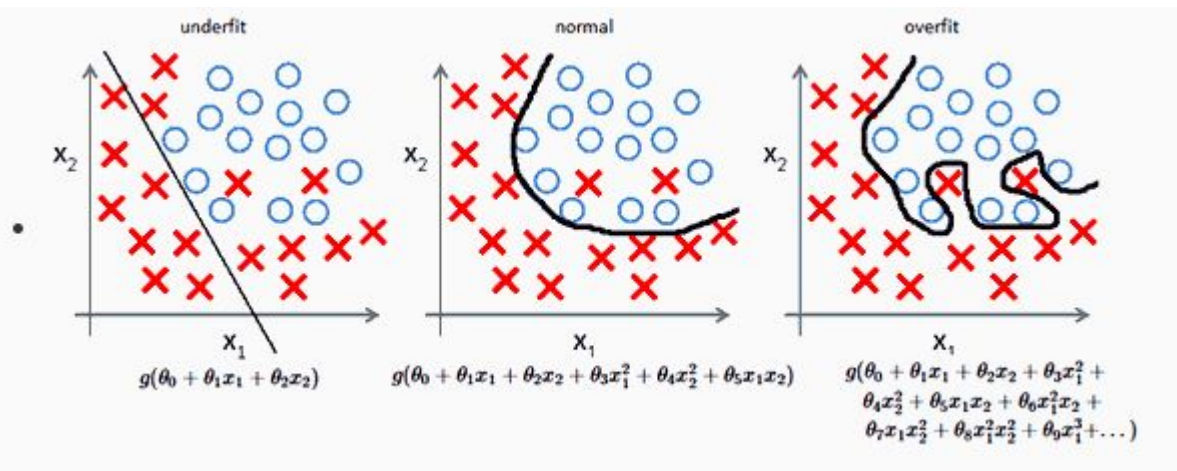


severely overfitting!



Overfitting

On Euclidean Space



Overfitting

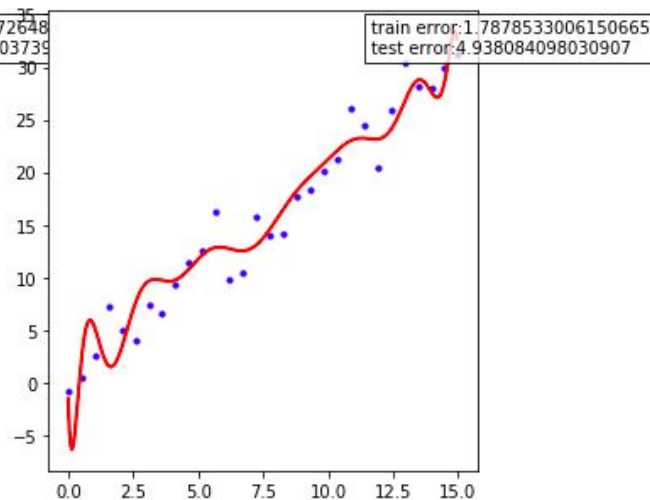
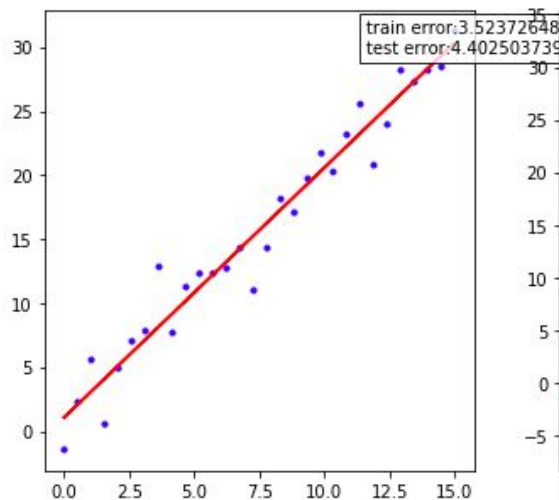
```
from model
lambda = 0.01

#generate test data
X_train = np.linspace(0, 15, 30)
y_train = np.random.normal(size=X_train.size)*np.sqrt(2)

#generate test data
X_test = np.linspace(0, 15, 20)
y_test = np.random.normal(size=X_test.size)*np.sqrt(2)

prop = np.polyfit(X_train, y_train, 1)
poly = np.polyval(prop, X_train)
overf = np.polyfit(X_train, y_train, 15)
overf = np.polyval(overf, X_test)

#plot
plt.figure(figsize=(10, 10))
plt.plot(X_train, y_train, 'b', label='train data')
plt.plot(X_test, y_test, 'b', label='test data')
plt.plot(X_train, poly, 'r', label='linear fit')
plt.plot(X_train, overf, 'r', label='overfit')
plt.legend()
plt.show()
```





The Problem of Overfitting

Too many complicated predictors in the model

- **Misleading statistics:** R-squared values, p-values, regression coefficients
- **Prediction error:** Not reflecting the overall population



Avoiding Overfitting: Regularization

Adding a complexity penalty to the loss function

1. L1 regularization (Lasso): sum of the weights

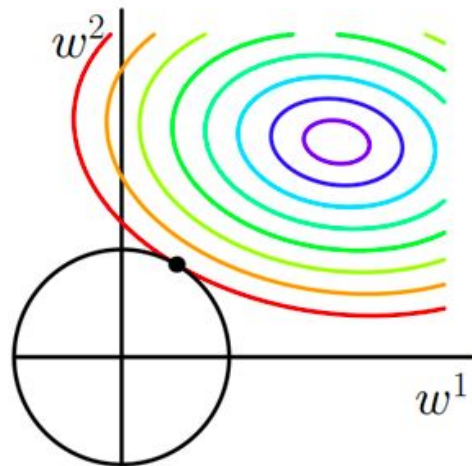
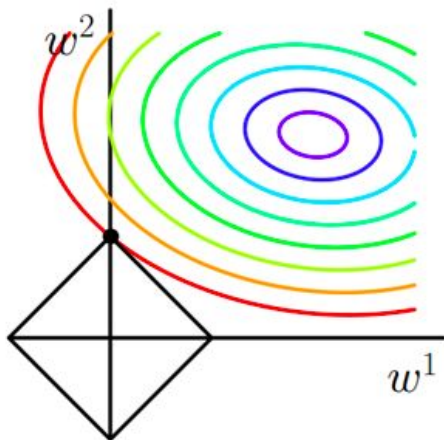
$$\min_{\beta} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1$$

2. L2 regularization (Ridge): sum of the square of the weights

$$\min_{\beta} \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$

Avoiding Overfitting: Regularization

On Euclidean Space



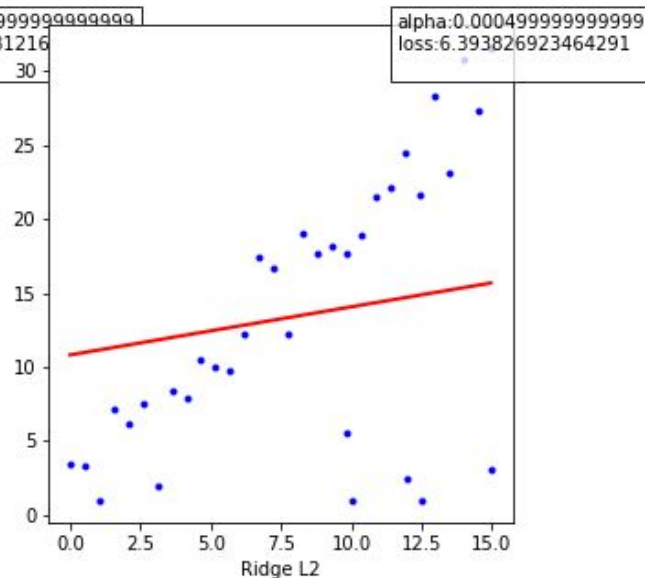
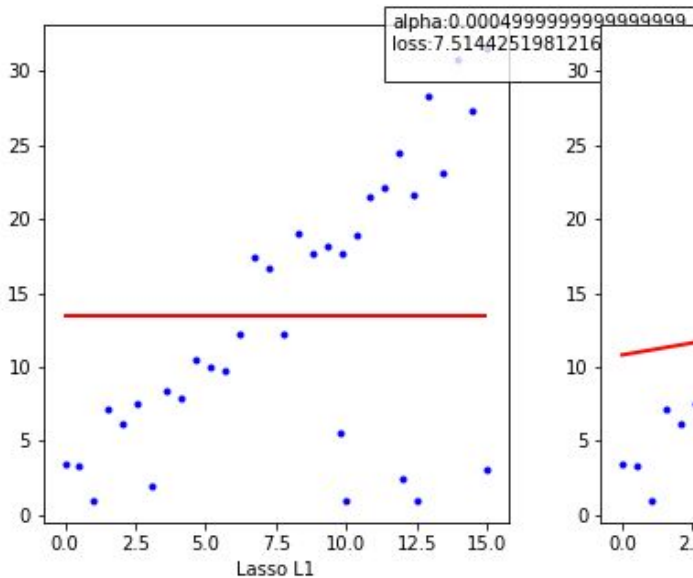


Avoiding Overfitting: Regularization

Bayesian Explanation (from MLE to MAP)

$$\underbrace{P(\theta \mid \mathcal{D})}_{\text{posterior}} \propto \underbrace{P(\mathcal{D} \mid \theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}$$

Time Step 0



```

#F1 classifiers
lassosco=list()
for alpha in alphas:
    thisLasso = Lasso(alpha=alpha)
    thisLasso.fit(x_train,y_train)
    lassoos.append(thisLasso.predict)

#F2 classifiers
ridgesco=list()
for alpha in alphas:
    thisRidge = Ridge(alpha=alpha)
    thisRidge.fit(x_train,y_train)
    ridges.append(thisRidge.predict)

#accuracies
lassosco_test_loss=(abs(func(x_test)-y_test.reshape(-1)).mean() for func in lassos)
ridgesco_test_loss=(abs(func(x_test)-y_test).mean() for func in ridges)

#generate animation
thisPlotter = Plotter('[[x_train, y_train]]*2')
plot1='func:lassosco, alpha:1/alphas, name: "Lasso 11", loss:lassosco_test_loss\'
plot2='func:ridgesco, alpha:1/alphas, name: "Ridge 11", loss:ridgesco_test_loss\'
plot3='func:lassosco, alpha:1/alphas, name: "Lasso 12", loss:lassosco_test_loss\'
plot4='func:ridgesco, alpha:1/alphas, name: "Ridge 12", loss:ridgesco_test_loss\'

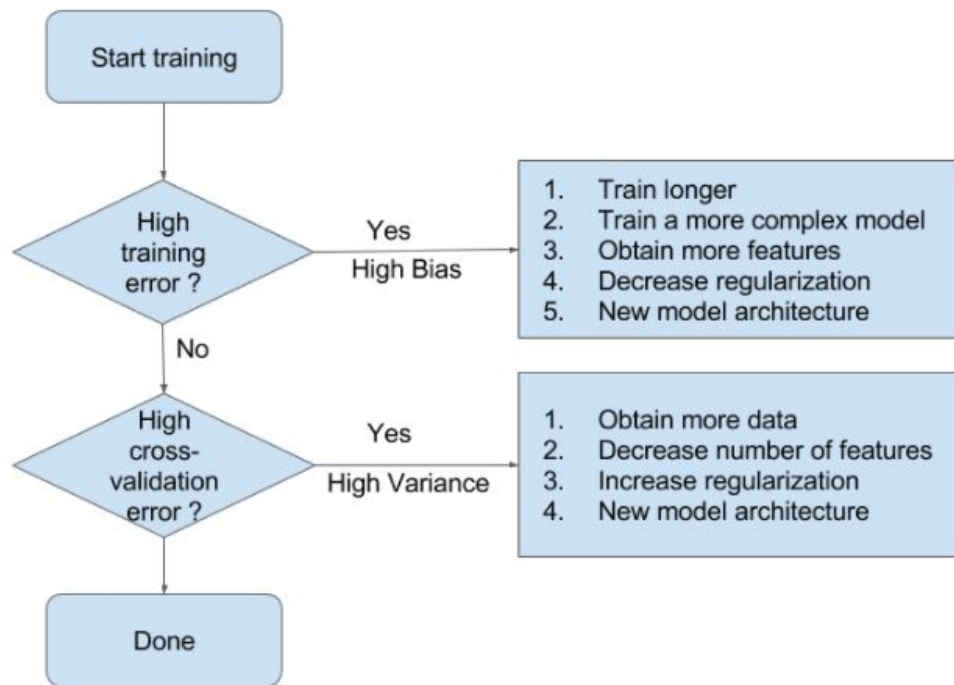
```



Hyper-Parameter Tuning and Evaluation

A Boston House Price Model
Enjoy!

Conclusion



Practical application of Linear regression

(a) Trend line

- Business analytics tool for forecasting , predicting And mostly monitoring tool used in various sector
- Technical analysis in share market is interpretation of trend line.



(b) crop yields on rainfall :

Yield is Dependent variable Rainfall is explanatory variable and by the use of data , prediction model are prepared using simple linear regression



(c) Economics

- To predict consumption/spending
- Fixed investment spending
- Inventory
- Export –import analysis
- Labor demand and supply
- Resource management
- CAPM





MANY MORE...

- SAT /GRE score for college admission
- Product price and sales
- Metal mining
- Impact of extraction on aquatic ecosystem
- Tobacco smoking v/s mortality

Q&A



THANK YOU !

