

Weekly Oxford Worldwide

DEPARTMENT FOR
CONTINUING
EDUCATION



PYTHON PROGRAMMING FOR DATA SCIENCE – PART 2 MASSIMILIANO IZZO & NICHOLAS DAY

LECTURE 3 SUPERVISED LEARNING: REGRESSION WITH SCIKIT-LEARN



Data preparation

- Feature engineering, where appropriate:
 - Discretize continuous features.
 - Decompose features (e.g., categorical, date/time, etc.).
 - Add promising transformations of features
 - Aggregate features into promising new features.
- Feature scaling:
 - Standardize or normalize features.
- Dimensionality Reduction (optional)
 - Principal Component Analysis (PCA)

Regression Performance Measure

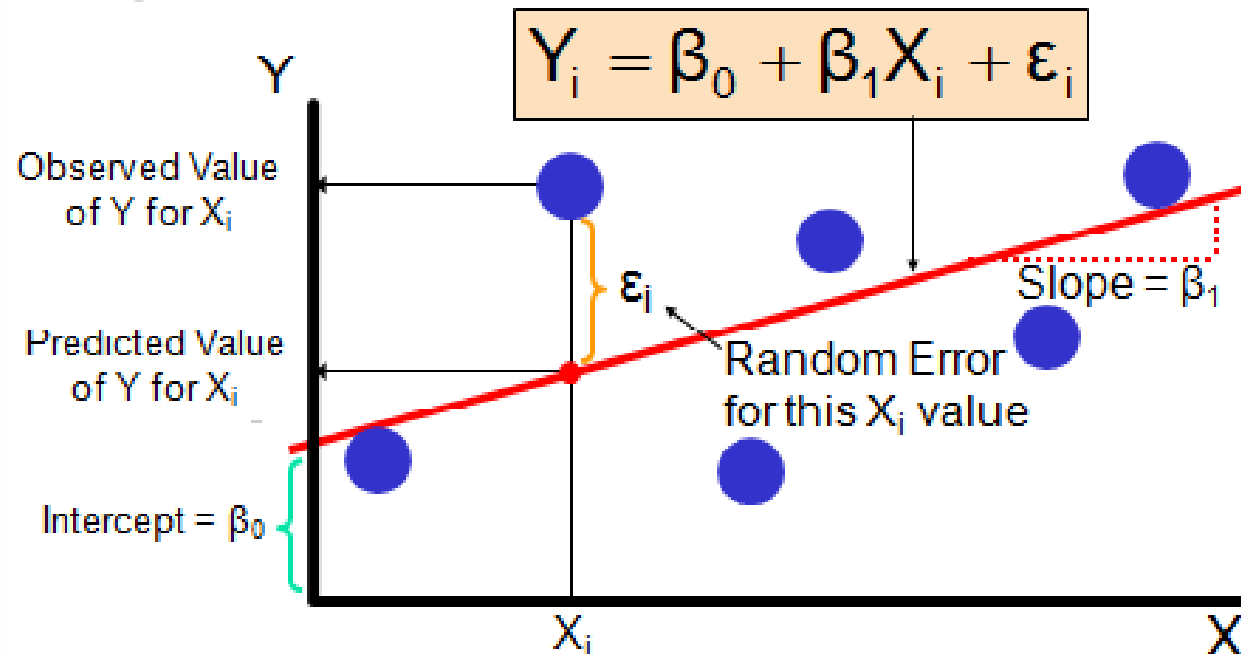
- $MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$
- $RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$
- $MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$
- $MedAE = median(|y_1 - \hat{y}_1|, |y_2 - \hat{y}_2|, \dots, |y_N - \hat{y}_N|)$
- $R^2 = 1 - \frac{MSE(model)}{MSE(baseline)} = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$

Choose some Algorithms

- Linear Regression: Ordinary Least Squares
 - Closed form solution (Normal Equation)
 - Gradient Descent
- Polynomial Regression
- Regularized Models
 - Ridge Regression
 - Lasso Regression
- Decision Trees Regression
- Something else (Support Vector Machines, Neural Networks...)
- Ensemble Models, Random Forests

Linear Regression

- Find the best linear model that fits our data
- This means finding two parameters: slope (β_1) and intercept (β_0)



Once trained, we
can use the model
to make predictions
=> machine
learning!!

Close form solution: Normal Equation

$$\hat{\beta} = \arg \min_{\beta} \|\underline{y} - \underline{X}\beta\|^2$$



$$\underline{X}^T \underline{X} \hat{\beta} = \underline{X}^T \underline{y}$$



$$\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Find the value of β that minimizes the squared sum of the estimation errors ϵ

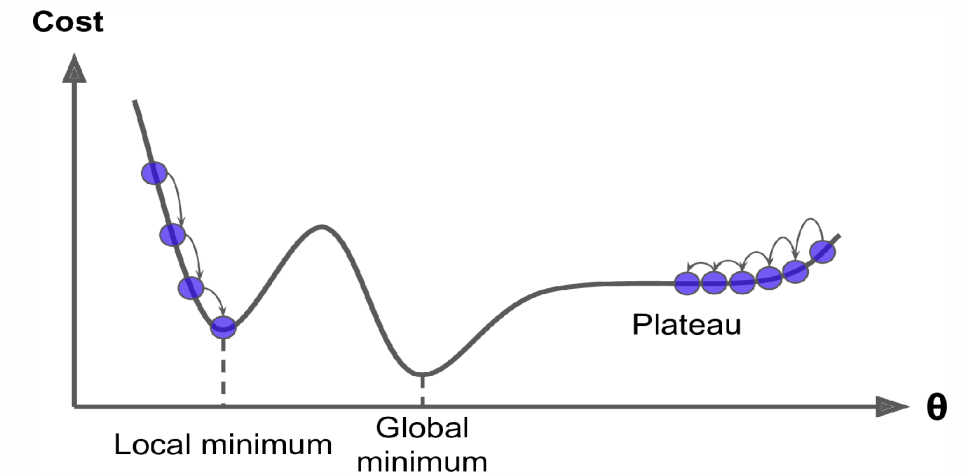
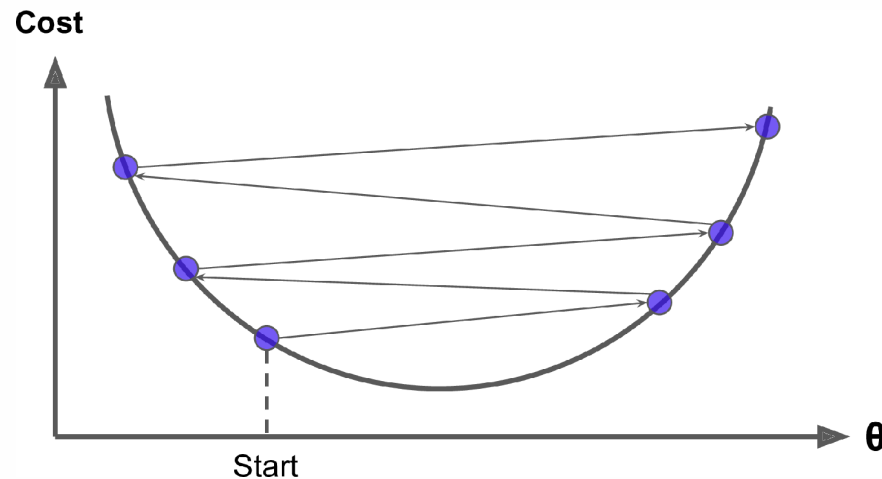
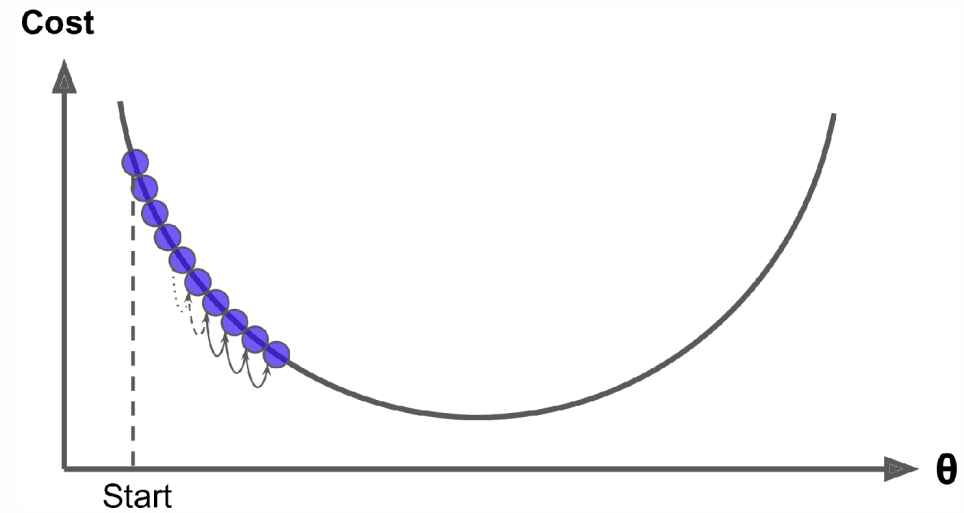
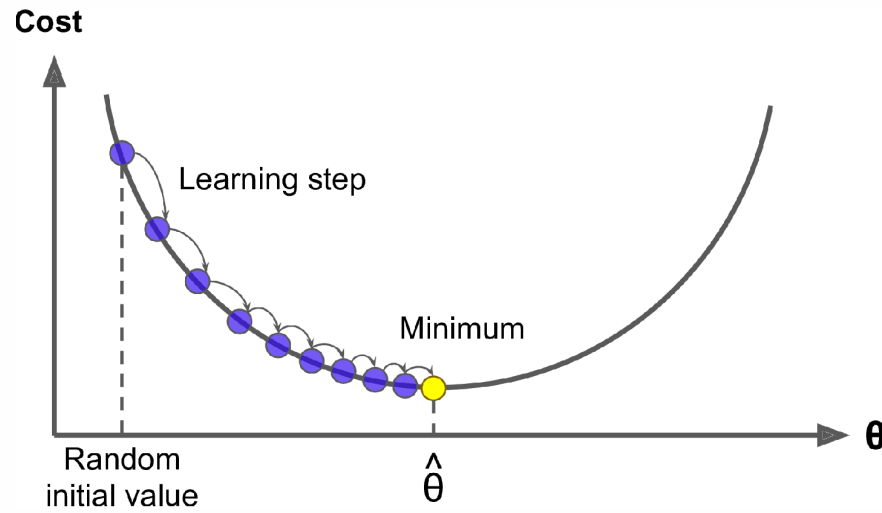
The issue here is the computational complexity of the explicit solution, especially the complexity of computing with respect to the number of features $(\underline{X}^T \underline{X})^{-1} \Rightarrow O(n^{2.4}) \div O(n^3)$

A different approach would be to use an optimisation algorithm to find the optimal solution

Gradient Descent

- Tweak the weights β iteratively in order to minimize a cost function.
- Measure the local gradient of the error function with respect to the weights β , and tweak β in the direction of descending gradient.
- Once the gradient equals zero, you have reached a minimum.

Gradient Descent



Batch Gradient Descent

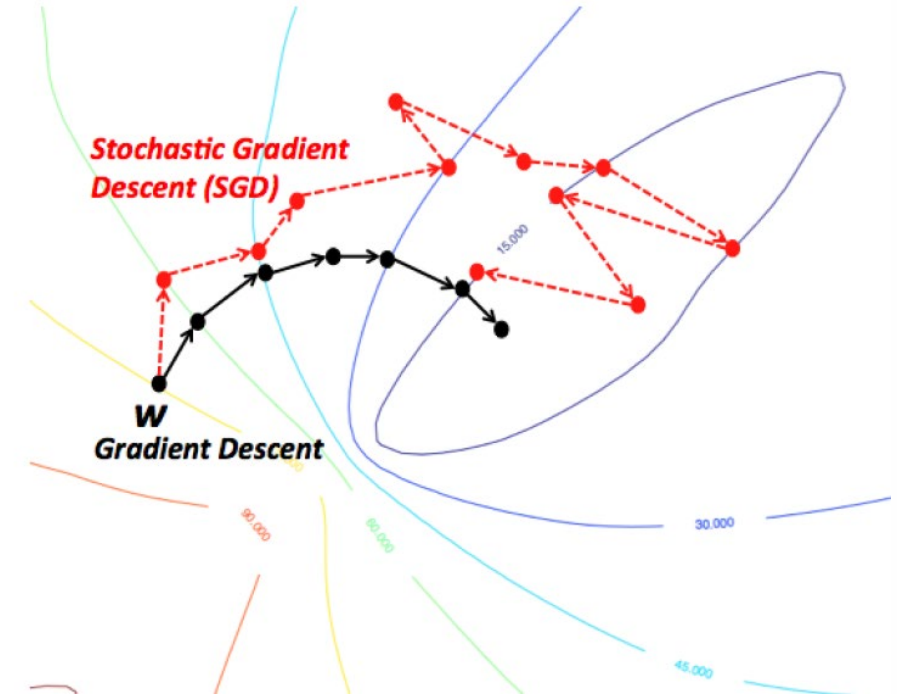
Let's consider the Mean Squared Error as our performance metric. The Gradient Formula for MSE is:

$$\nabla_{\theta} MSE(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} MSE(\theta) \\ \frac{\partial}{\partial \theta_2} MSE(\theta) \\ \dots \\ \dots \\ \frac{\partial}{\partial \theta_n} MSE(\theta) \end{pmatrix} = \frac{2}{m} X^T (X\theta - y)$$

Hence, the gradient descent at each iteration is: $\theta^{n+1} = \theta^n - \eta \nabla_{\theta} MSE(\theta)$

Stochastic Gradient Descent

- Note that this formula involves calculation over the full training set X at each Gradient Descent step. This is why the algorithm is called *Batch Gradient Descent*.
- The solution is to use *Stochastic Gradient Descent (SGD)*: pick a random instance in the training set at every step and computes the gradients based only on that single instance
- SGD: faster algorithm but slower to converge



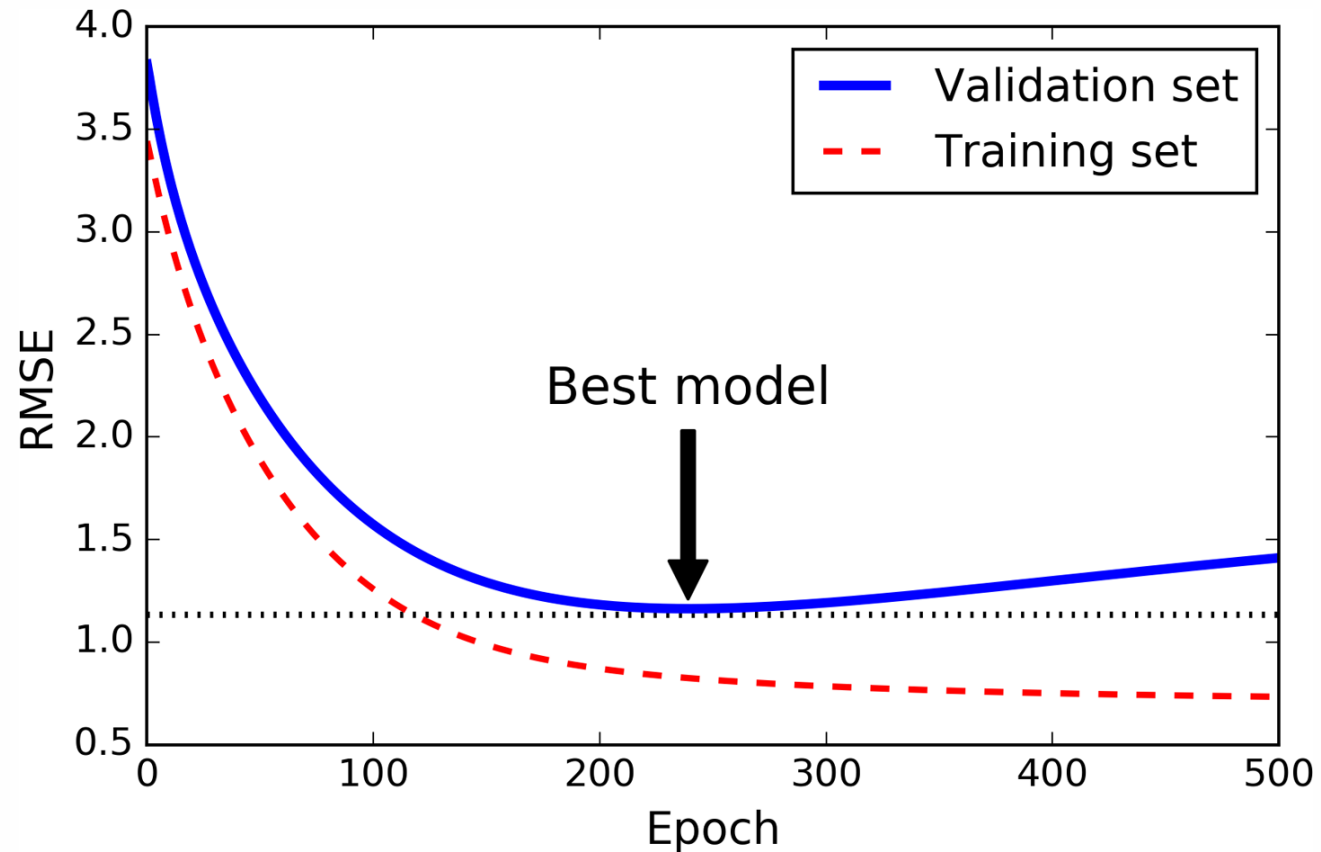
source: <https://wikidocs.net/3413>

Polynomial Regression and Regularisation

- My data is often more complex than a straight line (or a hyperplane) => need to add Polynomial Features to our fitting model
- P.R. is prone to overfitting => Regularization techniques are often needed
- Ridge Regression
- Lasso Regression
 - tends to eliminate the weights of the least important features
- ElasticNet (Ridge + Lasso)

Early Stopping

A very different way to regularize iterative learning algorithms such as Gradient Descent is to stop training as soon as the validation error reaches a minimum.



Cross-validation

