

# INFO251 – Applied Machine Learning

---

Lab 7  
Suraj R. Nair

# Announcements

- PS4 posted, due Monday March 12. Start early!
- Quiz 1 solutions in lecture tomorrow



# Today's Topics

1. Putting it all together: training an ML algorithm from scratch
  2. Common loss functions
  3. Practice
    - Bivariate OLS, squared error loss
    - Multivariate OLS, squared error loss
    - Multivariate OLS, squared error loss with Ridge regularization
  4. Cross validation for optimal regularization parameter
-

# Training ML Algorithms

Table 1: The three components of learning algorithms.

Representation	Evaluation	Optimization
Instances	Accuracy/Error rate	Combinatorial optimization
<i>K</i> -nearest neighbor	Precision and recall	Greedy search
Support vector machines	Squared error	Beam search
Hyperplanes	Likelihood	Branch-and-bound
Naive Bayes	Posterior probability	Continuous optimization
Logistic regression	Information gain	Unconstrained
Decision trees	K-L divergence	Gradient descent
Sets of rules	Cost/Utility	Conjugate gradient
Propositional rules	Margin	Quasi-Newton methods
Logic programs		Constrained
Neural networks		Linear programming
Graphical models		Quadratic programming
Bayesian networks		
Conditional random fields		

Domingos, 2016

# Training ML Algorithms (Linear Models)

1. Define a model
  2. Define a loss function
    - Add regularization to the loss function [**find optimal parameter**]
  3. Optimization
    - Gradient Descent [**batch size, step size / learning rate, stopping rule**]
-

# 1. Define a Model

- For today: **Linear regression models**, of the form  $y = ax + b$ 
    - Multivariate models:  $y = ax_1 + bx_2 + cx_3 + d$
    - Nonlinearities:  $y = ax^k + b$
    - Interaction terms:  $y = ax_1x_2 + b$
-

## 2. Define a Loss Function

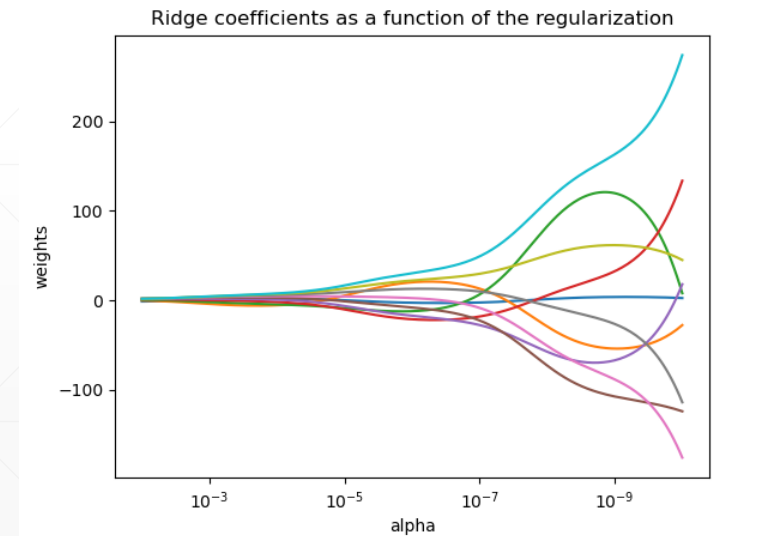
- Common loss functions for **regression**:
    - **Squared error loss**:  $J(y, \hat{y}) = (y - \hat{y})^2$
    - **Absolute error loss**:  $J(y, \hat{y}) = |y - \hat{y}|$
  - Common loss functions for **binary classification**:
    - **Logistic loss**:  $J(y, \hat{p}) = -(y \log(\hat{p}) + (1 - y) \log(1 - \hat{p}))$
    - **Hinge loss**:  $J(y, \hat{p}) = \max(0, 1 - \hat{p}y)$
  - Common loss functions for **multivariate classification**:
    - **Cross-entropy loss**:  $J(y, \hat{p}) = \sum_{c=1}^M y_c \log(\hat{p}_c)$
-

### 3. Optionally Add Regularization to the Loss

- **LASSO:**  $J(\theta) += \|\theta\|_1 = \sum_{j=1}^k |\theta_k|$



- **Ridge:**  $J(\theta) += \|\theta\|_2 = \sum_{j=1}^k \theta_k^2$





## 4. Gradient Descent

- Begin at a random point
  - Calculate the function value at the point and the gradient (partial derivatives)
  - Pick a new point, move in the direction of steepest descent. The size of the step is governed by the **learning rate**.
  - Repeat!
-

**Model:** Univariate least squares

**Cost:** Squared error

1. Define the model
  2. Define the loss function
  - ~~3. Optionally add regularization to the loss function~~
  4. Calculate partial derivatives
  5. Write pseudocode
-

**Model:** Multivariate least squares

**Cost:** Squared error

1. Define the model
  2. Define the loss function
  - ~~3. Optionally add regularization to the loss function~~
  4. Calculate partial derivatives
  5. Write pseudocode
-

**Model:** Multivariate least squares

**Cost:** Squared error + Ridge regularization

1. Define the model
  2. Define the loss function
  3. Optionally add regularization to the loss function
  4. Calculate partial derivatives
  5. Write pseudocode
-