

# Introduction to Machine Learning

Linear Regression

Varun Chandola

February 15, 2021

## Outline

## Contents

<b>1</b>	<b>Basics</b>	<b>1</b>
<b>2</b>	<b>Linear Regression</b>	<b>2</b>
2.1	Problem Formulation . . . . .	2
2.2	Learning Parameters . . . . .	4
2.3	Machine Learning as Optimization . . . . .	5
2.4	Convex Optimization . . . . .	6
2.5	Matrix Calculus Basics . . . . .	7
2.6	Evaluating Models . . . . .	7
2.7	Gradient Descent . . . . .	7
2.8	Issues with Gradient Descent . . . . .	10
2.9	Stochastic Gradient Descent . . . . .	11

## 1 Basics

### Basics

- Data - scalar ( $x$ ), vector ( $\mathbf{x}$ ), Matrix ( $\mathbf{X}$ )

### Scalars

- Numeric ( $x \in \mathbb{R}$ )
- Categorical (e.g.,  $x \in \{0, 1\}$ )
- Constants will be denoted as  $D$ ,  $M$ , etc.

### Vector

- Length of a vector  $\mathbf{x} \in \mathbb{R}^D$
- Vector *dot* product ( $\mathbf{x} \cdot \mathbf{y}$ )
- Norm of a vector ( $|\mathbf{x}|$ ,  $\|\mathbf{x}\|$ ,  $\|\mathbf{x}\|_p$ )

### Matrix

- Size of a matrix ( $\mathbf{X} \in \mathbb{R}^{M \times N}$ )
- Transpose of a matrix ( $\mathbf{X}^\top$ )
- Matrix product ( $\mathbf{XY}$ )
- A vector is a special matrix with only one column

$$\mathbf{x} \cdot \mathbf{y} \equiv \mathbf{x}^\top \mathbf{y}$$

## 2 Linear Regression

### 2.1 Problem Formulation

- There is one scalar **target** variable  $y$
- There is one vector **input** variable  $x$
- Inductive bias:

$$y = \mathbf{w}^\top \mathbf{x}$$

### Linear Regression Learning Task

Learn  $\mathbf{w}$  given training examples,  $\langle \mathbf{X}, \mathbf{y} \rangle$ .

The training data is denoted as  $\langle \mathbf{X}, \mathbf{y} \rangle$ , where  $\mathbf{X}$  is a  $N \times D$  data matrix consisting of  $N$  data examples such that each data example is a  $D$  dimensional vector.  $\mathbf{y}$  is a  $N \times 1$  vector consisting of corresponding target values for the examples in  $\mathbf{X}$ .

- Fitting a straight line to  $d$  dimensional data

$$y = \mathbf{w}^\top \mathbf{x}$$

$$y = \mathbf{w}^\top \mathbf{x} = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

- Will pass through origin
- Add intercept

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

- Equivalent to adding another column in  $\mathbf{X}$  of 1s.

### Incorporating Bias/Intercept

#### Explicit Bias

$$\begin{aligned} \mathbf{x} &\equiv \{x_1, x_2, \dots, x_d\} \\ \mathbf{w} &\equiv \{w_1, w_2, \dots, w_d\} \\ y &= w_0 + \mathbf{w}^\top \mathbf{x} \end{aligned}$$

#### Implicit Bias

$$\begin{aligned} \mathbf{x} &\equiv \{1, x_1, x_2, \dots, x_d\} \\ \mathbf{w} &\equiv \{w_0, w_1, w_2, \dots, w_d\} \\ y &= \mathbf{w}^\top \mathbf{x} \end{aligned}$$

## 2.2 Learning Parameters

- Minimize *squared loss*

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

- or,

$$J(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Make prediction  $(\mathbf{w}^\top \mathbf{x}_i)$  as close to the target  $(y_i)$  as possible
- Least squares estimate

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- *We will derive this expression in class.*

The derivation of the least squares estimate can be done by first converting the expression of the squared loss into matrix notation, i.e.,

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \\ &= \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \end{aligned}$$

To minimize the error, we first compute its derivative with respect to  $\mathbf{w}$ :

$$\frac{\partial LL(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Note that, we use the fact that  $(\mathbf{X}\mathbf{w})^\top \mathbf{y} = \mathbf{y}^\top \mathbf{X}\mathbf{w}$ , since both quantities are scalars and the transpose of a scalar is equal to itself. Continuing with the derivative:

$$\frac{\partial LL(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2} (2\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} - 2\mathbf{y}^\top \mathbf{X})$$

Setting the derivative to 0, we get:

$$\begin{aligned} 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} - 2\mathbf{y}^\top \mathbf{X} &= 0 \\ \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} &= \mathbf{y}^\top \mathbf{X} \\ (\mathbf{X}^\top \mathbf{X})^\top \mathbf{w} &= \mathbf{X}^\top \mathbf{y} \text{ (Taking transpose both sides)} \\ (\mathbf{X}^\top \mathbf{X}) \mathbf{w} &= \mathbf{X}^\top \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

## 2.3 Machine Learning as Optimization

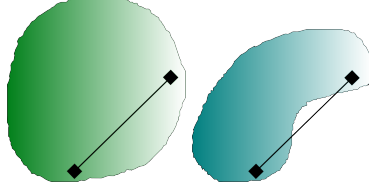
At this point, we move away from the situation where a perfect solution exists, and the learning task is to reach the perfect solution. Instead, we focus on finding the *best possible* solution which optimizes certain criterion.

- Learning is optimization
- Faster optimization methods for faster learning
- Let  $\mathbf{w} \in \mathbb{R}^d$  and  $S \subset \mathbb{R}^d$  and  $f_0(\mathbf{w}), f_1(\mathbf{w}), \dots, f_m(\mathbf{w})$  be real-valued functions.
- Standard optimization formulation is:

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && f_0(\mathbf{w}) \\ & \text{subject to} && f_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, m. \end{aligned}$$

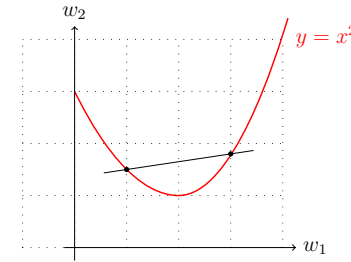
- Methods for **general optimization problems**
  - Simulated annealing, genetic algorithms
- Exploiting *structure* in the optimization problem
  - **Convexity**, Lipschitz continuity, smoothness

### Convex Sets



### Convex Functions

<sup>1</sup>Adapted from [http://ttic.uchicago.edu/~gregory/courses/ml2012/lectures/tutorial\\_optimization.pdf](http://ttic.uchicago.edu/~gregory/courses/ml2012/lectures/tutorial_optimization.pdf). Also see, <http://www.stanford.edu/~boyd/cvxbook/> and [http://scipy-lectures.github.io/advanced/mathematical\\_optimization/](http://scipy-lectures.github.io/advanced/mathematical_optimization/).



*Convexity* is a property of certain functions which can be exploited by optimization algorithms. The idea of convexity can be understood by first considering *convex sets*. A convex set is a set of points in a coordinate space such that every point on the line segment joining any two points in the set are also within the set. Mathematically, this can be written as:

$$\mathbf{w}_1, \mathbf{w}_2 \in S \Rightarrow \lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2 \in S$$

where  $\lambda \in [0, 1]$ . A *convex function* is defined as follows:

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex function if the domain of  $f$  is a convex set and for all  $\lambda \in [0, 1]$ :

$$f(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2) \leq \lambda f(\mathbf{w}_1) + (1 - \lambda) f(\mathbf{w}_2)$$

Some examples of convex functions are:

- Affine functions:  $\mathbf{w}^\top \mathbf{x} + b$
- $\|\mathbf{w}\|_p$  for  $p \geq 1$
- Logistic loss:  $\log(1 + e^{-y\mathbf{w}^\top \mathbf{x}})$

## 2.4 Convex Optimization

- Optimality Criterion

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimize}} && f_0(\mathbf{w}) \\ & \text{subject to} && f_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, m. \end{aligned}$$

where all  $f_i(\mathbf{w})$  are **convex functions**.

- $\mathbf{w}_0$  is feasible if  $\mathbf{w}_0 \in \text{Dom } f_0$  and all constraints are satisfied
- A feasible  $\mathbf{w}^*$  is optimal if  $f_0(\mathbf{w}^*) \leq f_0(\mathbf{w})$  for all  $\mathbf{w}$  satisfying the constraints

## 2.5 Matrix Calculus Basics

$$\frac{\partial \mathbf{a}^\top \mathbf{b}}{\partial \mathbf{a}} = \frac{\partial \mathbf{b}^\top \mathbf{a}}{\partial \mathbf{a}} = \mathbf{b}$$

$$\frac{\partial \mathbf{a}^\top \mathbf{M} \mathbf{a}}{\partial \mathbf{a}} = 2\mathbf{M} \mathbf{a}$$

where  $\mathbf{M}$  is a symmetric matrix.

## 2.6 Evaluating Models

### Evaluating Linear Regression Model

- How do we know that a model is good?
- What is a good evaluation/performance metric?

#### Root Mean Squared Error

$$RMSE = \sqrt{\sum_{i=1}^N (y^i - \hat{y}^i)^2}$$

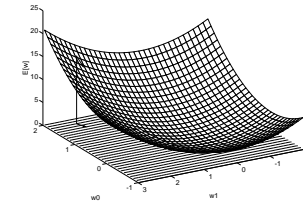
where  $\hat{y}_i$  is the prediction for the  $i^{th}$  instance.

- What data to evaluate this on?
  - Training data?
  - Test data (**generalization error**)

## 2.7 Gradient Descent

- Denotes the direction of steepest ascent

$$\nabla J(\mathbf{w}) = \begin{bmatrix} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_d} \end{bmatrix}$$



A small step in the weight space from  $\mathbf{w}$  to  $\mathbf{w} + \delta \mathbf{w}$  changes the objective (or error) function. This change is maximum if  $\delta \mathbf{w}$  is along the direction of the gradient at  $\mathbf{w}$  and is given by:

$$\delta J \simeq \delta \mathbf{w}^\top \nabla J(\mathbf{w})$$

Since  $J(\mathbf{w})$  is a smooth continuous function of  $\mathbf{w}$ , the extreme values of  $J$  will occur at the location in the input space ( $\mathbf{w}$ ) where the gradient of the error function vanishes, such that:

$$\nabla J(\mathbf{w}) = 0$$

The vanishing points can be further analyzed to identify them as saddle, minima, or maxima points.

One can also derive the local approximations done by first order and second order methods using the Taylor expansion of  $J(\mathbf{w})$  around some point

$\mathbf{w}'$  in the weight space.

$$J(\mathbf{w}') \simeq J(\mathbf{w}) + (\mathbf{w}' - \mathbf{w})^\top \nabla + \frac{1}{2}(\mathbf{w}' - \mathbf{w})^\top \mathbf{H}(\mathbf{w}' - \mathbf{w})$$

For first order optimization methods, we ignore the second order derivative (denoted by  $\mathbf{H}$  or the *Hessian*). It is easy to see that for  $\mathbf{w}$  to be the local minimum,  $J(\mathbf{w}) - J(\mathbf{w}') \leq 0$ ,  $\forall \mathbf{w}'$  in the vicinity of  $\mathbf{w}$ . Since we can choose any arbitrary  $\mathbf{w}'$ , it means that every component of the gradient  $\nabla$  must be zero.

1. Set derivative to 0

$$\nabla J(\mathbf{w}) = 0$$

2. Check second derivative for minima or maxima or saddle point

1. Start from any point in variable space

2. Move along the direction of the steepest descent (or ascent)

- By how much?
- A learning rate ( $\eta$ )
- What is the direction of steepest descent?
  - Gradient of  $J$  at  $\mathbf{w}$

Gradient descent is a first-order optimization method for convex optimization problems. It is analogous to “hill-climbing” where the gradient indicates the direction of steepest ascent in the local sense.

#### Training Rule for Gradient Descent

$$\mathbf{w} = \mathbf{w} - \eta \nabla J(\mathbf{w})$$

For each weight component:

$$w_j = w_j - \eta \frac{\partial J}{\partial w_j}$$

The key operation in the above update step is the calculation of each partial derivative. This can be computed for perceptron error function as

follows:

$$\begin{aligned} \frac{\partial J}{\partial w_j} &= \frac{\partial}{\partial w_j} \frac{1}{2} \sum_i (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \\ &= \frac{1}{2} \sum_i \frac{\partial}{\partial w_j} (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \\ &= \frac{1}{2} \sum_i 2(y_i - \mathbf{w}^\top \mathbf{x}_i) \frac{\partial}{\partial w_j} (y_i - \mathbf{w}^\top \mathbf{x}_i) \\ &= \sum_i (y_i - \mathbf{w}^\top \mathbf{x}_i) (-x_{ij}) \end{aligned}$$

where  $x_{ij}$  denotes the  $j^{\text{th}}$  attribute value for the  $i^{\text{th}}$  training example. The final weight update rule becomes:

$$w_j = w_j + \eta \sum_i (y_i - \mathbf{w}^\top \mathbf{x}_i) x_{ij}$$

- Error surface contains only one global minimum
- Algorithm *will* converge
  - Examples need not be linearly separable
- $\eta$  should be *small enough*
- Impact of too large  $\eta$ ?
- Too small  $\eta$ ?

If the learning rate is set very large, the algorithm runs the risk of overshooting the target minima. For very small values, the algorithm will converge very slowly. Often,  $\eta$  is set to a moderately high value and reduced after each iteration.

## 2.8 Issues with Gradient Descent

- Slow convergence
- Stuck in local minima

One should note that the second issue will not arise in the case of Perceptron training as the error surface has only one global minima. But for general setting, including multi-layer perceptrons, this is a typical issue.

More efficient algorithms exist for batch optimization, including *Conjugate Gradient Descent* and other *quasi*-Newton methods. Another approach is to consider training examples in an online or incremental fashion, resulting in an online algorithm called **Stochastic Gradient Descent** [1], which will be discussed next.

## 2.9 Stochastic Gradient Descent

- Update weights after every training example.
- For sufficiently small  $\eta$ , closely approximates Gradient Descent.

Gradient Descent	Stochastic Gradient Descent
Weights updated after summing error over all examples	Weights updated after examining each example
More computations per weight update step	Significantly lesser computations
Risk of local minima	Avoids local minima

- Minimize the squared loss using *Gradient Descent*

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

- Why?

The analytical approach discussed earlier involves a matrix inversion  $((\mathbf{X}^\top \mathbf{X})^{-1})$  which is a  $(D+1) \times (D+1)$  matrix. Alternatively, one could solve a system of equations. When  $D$  is large, this inversion can be computationally expensive ( $O * D^3$ ) for standard matrix inversion. Moreover, often, the linear system might have singularities and inversion or solving the system of equations might yield numerically unstable results.

To compute the gradient update rule one can differentiate the error with respect to each entry of  $\mathbf{w}$ .

$$\begin{aligned} \frac{\partial J(\mathbf{w})}{\partial w_j} &= \frac{1}{2} \frac{\partial}{\partial w_j} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \\ &= \sum_{i=1}^N (\mathbf{w}^\top \mathbf{x}_i - y_i) x_{ij} \end{aligned}$$

Using the above result, one can perform repeated updates of the weights:

$$w_j := w_j - \eta \frac{\partial J(\mathbf{w})}{\partial w_j}$$

## References

## References

- [1] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Comput.*, 1(4):541–551, Dec. 1989.