

# Machine Learning for Intelligent Systems

Lecture 12: Stochastic Gradient Descent

Reading: UML 14

Instructors: Nika Haghtalab (this time) and Thorsten Joachims

1

## Regularized Linear Models

Many learning problems can be written as the following optimization on the sample set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .

$$\min_{\vec{w}, b} \underbrace{R(\vec{w})}_{\text{Regularizer}} + C \frac{1}{n} \sum_{i=1}^n \underbrace{L(\vec{w} \cdot \vec{x}_i + b, y_i)}_{\text{Loss of each instance } (\vec{x}_i, y_i)}$$

1. Primal SVM:  $\min_{\vec{w}, b} \frac{1}{2} \vec{w} \cdot \vec{w} + C \frac{1}{n} \sum_{i=1}^n \max(1 - y_i (\vec{w} \cdot \vec{x}_i + b), 0)$
2. Reg. Logistic Regression:  $\min_{\vec{w}, b} \frac{1}{2} \vec{w} \cdot \vec{w} + C \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i (\vec{w} \cdot \vec{x}_i + b)})$
3. Ridge Regression:  $\min_{\vec{w}, b} \frac{1}{2} \vec{w} \cdot \vec{w} + C \frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i + b - y_i)^2$

2

## Loss Functions

$$\min_{\vec{w}, b} R(\vec{w}) + C \frac{1}{n} \sum_{i=1}^n L(\vec{w} \cdot \vec{x}_i + b, y_i)$$

| Loss Function $L(y, y_i)$                 | Algorithm           |
|---|---------------------|
| <b>Hinge loss:</b> $\max(1 - y y_i, 0)$   | SVM                 |
| <b>Log loss:</b> $\log(1 + \exp(-y y_i))$ | Logistic Regression |
| <b>Exponential loss:</b> $\exp(-y y_i)$   | Boosting            |
| <b>0-1 Loss:</b> $\mathbf{1}(y \neq y_i)$ | Classification loss |

Non-Convex

Convex

3

## Regularizers

$$\min_{\vec{w}, b} R(\vec{w}) + C \frac{1}{n} \sum_{i=1}^n L(\vec{w} \cdot \vec{x}_i + b, y_i)$$

| Regularizer $R(\vec{w})$                                     | Properties              |
|--|-------------------------|
| $\ell_2$ regularization: $\frac{1}{2} \vec{w} \cdot \vec{w}$ | Convex                  |
| $\ell_1$ regularization: $\ \vec{w}\ _1$                     | Convex, sparse          |
| $\ell_p$ for $0 \leq p < 1$                                  | Non-convex, very sparse |

4

## Optimizing Regularized Lin. Models

Many learning problems can be written as the following optimization on the sample set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .

$$\min_{\vec{w}} \underbrace{R(\vec{w}) + C \frac{1}{n} \sum_{i=1}^n L(\vec{w} \cdot \vec{x}_i, y_i)}_{\mathcal{L}_S(\vec{w})}$$

Formal guarantees for when  $\mathcal{L}_S(\vec{w})$  is convex in  $\vec{w}$ . But these methods are widely used for non-convex optimization as well.

5

## Gradient Descent (GD)

For finding the minimum of a convex function:  $\min_{\vec{w}} \mathcal{L}(\vec{w})$

**Gradient**

$$\nabla \mathcal{L}(\vec{w}) = \left( \frac{\partial \mathcal{L}(\vec{w})}{\partial w_1}, \dots, \frac{\partial \mathcal{L}(\vec{w})}{\partial w_d} \right)$$

**Gradient Descent**

**Input:** A function  $\mathcal{L}$ , number of time steps  $T$ , step size  $\eta_t$

**Initialize**  $\vec{w}^{(0)} = (0, \dots, 0)$

**For**  $t = 1, 2, 3, \dots, T$

$$\vec{w}^{(t+1)} \leftarrow \vec{w}^{(t)} - \eta_t \nabla \mathcal{L}(\vec{w}^{(t)})$$

**Output**  $\vec{w}^{(T)}$

6

### Learning with Gradient Descent

Consider the following optimization on  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .

$$\min_{\bar{w}} \mathcal{L}_S(\bar{w}) = \min_{\bar{w}} R(\bar{w}) + C \frac{1}{n} \sum_{i=1}^n L(\bar{w} \cdot \bar{x}_i, y_i)$$

**Gradient Descent for Learning**

**Input:** Step sizes  $\eta_t$ , time  $T$ , and samples  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$   
**Initialize**  $\bar{w}^{(0)} = (0, \dots, 0)$   
**For**  $t = 1, 2, 3, \dots, T$

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t \frac{C}{n} \sum_{i=1}^n \nabla L(\bar{w}^{(t)} \cdot \bar{x}_i, y_i)$$

**Output**  $\bar{w}^{(T)}$

7

### Example: GD for SVM

**Example:** Consider the SVM primal form as written in a regularized linear model (homogenous).

$$\min_{\bar{w}} \frac{1}{2} \bar{w} \cdot \bar{w} + C \frac{1}{n} \sum_{i=1}^n \max(1 - y_i(\bar{w} \cdot \bar{x}_i), 0)$$

Gradient of

$$\nabla R(\bar{w}) = \bar{w} \quad \nabla L_S(\bar{w}) = \frac{C}{n} \sum_{i=1}^n \text{???}$$

$\max(1 - y_i(\bar{w} \cdot \bar{x}_i), 0)$   
 0 if  $y_i(\bar{w} \cdot \bar{x}_i) > 1$ ,       $1 - y_i(\bar{w} \cdot \bar{x}_i)$  if  $y_i(\bar{w} \cdot \bar{x}_i) \leq 1$   
Gradient = 0                      Gradient =  $-y_i x_i$

**GD update:**  $\bar{w}^{(t+1)} \leftarrow (1 - \eta_t) \bar{w}^{(t)} + \frac{\eta_t C}{n} \sum_{i=1}^n y_i x_i \mathbf{1}(y_i(\bar{w}^{(t)} \cdot \bar{x}_i) \leq 1)$

8

### Gradient Descent for Large Scale ML

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t \frac{C}{n} \sum_{i=1}^n \nabla L(\bar{w}^{(t)} \cdot \bar{x}_i, y_i)$$

Challenges?

- Large data set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  for very large  $n$ .
- High dimensional data sets  $x_i \in \mathbb{R}^d$  for very large  $d$ .

9

### Using fewer samples for the update

Each time step use fewer samples for update

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t \frac{C}{n} \sum_{i=1}^n \nabla L(\bar{w}^{(t)} \cdot \bar{x}_i, y_i)$$

↓

Take a random  $(\bar{x}_{(t)}, y_{(t)}) \sim S$ .

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t C \nabla L(\bar{w}^{(t)} \cdot \bar{x}_{(t)}, y_{(t)})$$

10

### Stochastic Gradient Descent (SGD)

**Gradient Descent for Learning**

**Input:** Function  $\mathcal{L}$ , step sizes  $\eta_t$ , time  $T$ , samples  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$   
**Initialize**  $\bar{w}^{(0)} = (0, \dots, 0)$   
**For**  $t = 1, 2, 3, \dots, T$   
 Take a random sample of  $(\bar{x}, y) \sim S$   

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t C \nabla L(\bar{w}^{(t)} \cdot \bar{x}, y)$$

**Output**  $\bar{w}^{(T)}$ .      %% Or the average of  $\bar{w}^{(1)}, \dots, \bar{w}^{(T)}$

Stochastic Gradient Descent:

— Each iteration  
 — Average up to now

11

Added from the whiteboard

### Example: SGD for SVM

**Example:** Consider the SVM primal form as written in a regularized linear model (homogenous).

$$\min_{\bar{w}} \frac{1}{2} \bar{w} \cdot \bar{w} + C \frac{1}{n} \sum_{i=1}^n \max(1 - y_i(\bar{w} \cdot \bar{x}_i), 0)$$

**SGD update:** Take  $(\bar{x}_i, y_i) \sim S$

$$\bar{w}^{(t+1)} \leftarrow (1 - \eta_t) \bar{w}^{(t)} + \eta_t C y_i \bar{x}_i \mathbf{1}(y_i(\bar{w}^{(t)} \cdot \bar{x}_i) \leq 1)$$

**Equivalently:**  
 Take  $(\bar{x}_i, y_i) \sim S$   
 If  $y_i(\bar{w}^{(t)} \cdot \bar{x}_i) \leq 1$  then  $\bar{w}^{(t+1)} \leftarrow (1 - \eta_t) \bar{w}^{(t)} + \eta_t C y_i \bar{x}_i$   
 Else  $\bar{w}^{(t+1)} \leftarrow (1 - \eta_t) \bar{w}^{(t)}$ .

12

### Effects of Step Size

Small step size

Medium step size

Larger step size

Smaller step size: More similar to Gradient Descent, less stochastic improvement, less uncertainty  
 Bigger step size: More stochastic improvement, more uncertainty.

13

### What makes SGD work

Gradient Descent:

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t \underbrace{C \frac{1}{n} \sum_{i=1}^n \nabla L(\bar{w}^{(t)} \cdot \vec{x}_i, y_i)}_{E_{(\vec{x}, y)}[\nabla L(\bar{w}^{(t)} \cdot \vec{x}, y)]}$$

Stochastic Gradient Descent:

- Uses an “unbiased” estimator for the total gradient.
- Step size helps control the variance.

Noisy Estimates

Reduce Computation

Improve Generalization

14

### Mini-Batch

Go between deterministic Gradient Descent and Stochastic Gradient Descent, take between  $n$  and  $1$  instances.

At each time: Take a random subset  $S_t$  of instance

$$\bar{w}^{(t+1)} \leftarrow \bar{w}^{(t)} - \eta_t \nabla R(\bar{w}^{(t)}) - \eta_t C \frac{1}{|S_t|} \sum_{(x,y) \in S_t} \nabla L(\bar{w}^{(t)} \cdot \vec{x}, y)$$

— Batch gradient descent

— Mini-batch gradient Descent

— Stochastic gradient descent

15

### Practical Challenges

Randomly choosing an instance or mini-batch

- In practice: Shuffle and choose without replacement
- Theory: i.i.d, with replacement

Beyond the linear  $\bar{w} \cdot \vec{x}$

- Convex Versus Non-Convex

Setting the step size and mini-batch size?

- Parallel computation.
- Generalization?
- Convex versus Non-convex

16