Machine Learning (CSE 446): Learning as Minimizing Loss: Regularization and Gradient Descent

> Sham M Kakade © 2018

University of Washington cse446-staff@cs.washington.edu

イロト イロト イヨト イヨト 二日

1/12

### Announcements

- Assignment 2 due tomo.
- ► Midterm: Weds, Feb 7th.
- ► Qz section: review

#### Today:

Regularization and Optimization!

# Review

### Relax!

► The mis-classification optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} \llbracket y_n(\mathbf{w} \cdot \mathbf{x}_n) \le 0 \rrbracket$$

• Instead, use loss function  $\ell(y_n, \mathbf{w} \cdot \mathbf{x})$  and solve a**relaxation**:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n)$$

### Relax!

► The mis-classification optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} \llbracket y_n(\mathbf{w} \cdot \mathbf{x}_n) \le 0 \rrbracket$$

• Instead, use loss function  $\ell(y_n, \mathbf{w} \cdot \mathbf{x})$  and solve a**relaxation**:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n)$$

- ► What do we want? 🧲
- How do we get it? speed? accuracy?

# Some loss functions:

► The square loss:

$$\ell(y, \mathbf{w} \cdot \mathbf{x}) = (y - \mathbf{w} \cdot \mathbf{x})^2$$

► The logistic loss:

$$-\ell^{\text{logistic}}(y, \mathbf{w} \cdot \mathbf{x}) = \log\left(1 + \exp(-y\mathbf{w} \cdot \mathbf{x})\right).$$

- ► They both "upper bound" the mistake rate.
- Instead:
  - $\blacktriangleright$  Instead, we let's care about "regression" where y is real valued.
  - What if we have multiple classes? (not just binary classification?)

1

9≈ W.7

Least squares: let's minimize it!

► The optimization problem:

where Y is an h-vector and X is our  $h \times d$  data matrix.

• The solution is the **least squares estimator**:

$$\mathbf{w}^{\text{least squares}} = (X^{\top}X)^{-1}X^{\top}Y$$

4 ロ ト 4 日 ト 4 目 ト 4 目 ト 目 の Q (\*
4/12)

Matrix calculus proof: scratch space  $\|X w - Y\|^{2} = (X w - Y)^{T} (X w - Y) = w^{T} X^{T} x w - X^{T} X w + Y^{T} Y$  $\frac{\partial}{\partial z} \left( \right) = \chi X' X \omega - \chi X' Y$ want  $(X^T X) w = X^T Y$ 

Y = Xn

<ロ> < (回) < (u) < (u) </p>

Matrix calculus proof: scratch space

<ロト <部ト < Eト < Eト 差 の < で 5/12 Let's remember our linear system solving!

 $2 \psi_{1} + 5 \psi_{2} = 8$   $3 \psi_{1} + \frac{10 \psi_{2}}{2} = 11$   $\left( \begin{array}{c} 2 \\ 3 \\ 3 \\ 3 \\ 3 \\ 10 \end{array} \right) \left( \begin{array}{c} \psi_{1} \\ \psi_{2} \end{array} \right) = \left( \begin{array}{c} 8 \\ 11 \end{array} \right)$ 

# Today

Least squares: What could go wrong?!

► The optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} (y_n - \mathbf{w} \cdot \mathbf{x}_n)^2 = \\\min_{\mathbf{w}} \|Y - X\mathbf{w}\|^2$$

where Y is an n-vector and X is our  $n \times d$  data matrix.

• The solution is the **least squares estimator**:

$$\mathbf{w}^{\text{least squares}} = (X^{\top}X)^{-1}X^{\top}Y$$

Least squares: What could go wrong?!

► The optimization problem:

$$\begin{split} \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} (y_n - \mathbf{w} \cdot \mathbf{x}_n)^2 = \\ \min_{\mathbf{w}} \|Y - X\mathbf{w}\|^2 \end{split}$$
 where Y is an A-vector and X is our  $\mathbf{w} \times d$  data matrix.

• The solution is the **least squares estimator**:

$$\mathbf{w}^{\text{least squares}} = (X^{\top}X)^{-1}X^{\top}Y$$

What if d is bigger than h? Even if not?

# What could go wrong?

Suppose d > h: (XTX) - not invertible

solve e.g. 5w,+3wz=1) (under mostraned system)

What about n > d?

<ロ> < (回) < (0) </p>
8/12

# What could go wrong?

Suppose d > h:

What about p > d?

- What happens if features are very correlated?
  - (e.g. 'rows/columns in our matrix are **co-linear**.)

exactly

イロト 不得下 イヨト イヨト 二日

8/12

linear system solving: scratch space

$$2W_1 + 3W_2 = 11$$
  
 $3W_1 + 3.0003W_2 = 1081$ 

A fix: Regularization

Regularize the optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} (y_n - \mathbf{w} \cdot \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|^2 = \sum_{n=1}^{\infty} \|Y - X^{\top} \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

► This particular case: "Ridge" Regression, Tikhonov regularization

• The solution is the **least squares estimator**:

$$\mathbf{w}^{\text{least squares}} = \left(\frac{1}{N}X^{\top}X + \mathbf{II}\right)^{-1} \left(\frac{1}{N}X^{\top}Y\right)$$

2 regu

# The "general" approach

► The **regularized** optimization problem:

$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, \mathbf{w} \cdot \mathbf{x}_n) + R(\mathbf{w})$$

▶ Penalty some w more than others. Example:  $R(w) = ||w||^2$ 

How do we find a solution quickly?

# Remember: convexity



<ロト < 部 > < 言 > < 言 > こ > < 言 > こ ? へ (~ 10 / 12

# Gradient Descent



### Gradient Descent

Data: function  $F : \mathbb{R}^d \to \mathbb{R}$ , number of iterations K, step sizes  $\langle \eta^{(1)}, \ldots, \eta^{(K)} \rangle$ Result:  $\mathbf{z} \in \mathbb{R}^d$ initialize:  $\mathbf{z}^{(0)} = \mathbf{0}$ ; for  $k \in \{1, \ldots, K\}$  do  $| \mathbf{z}^{(k)} = \mathbf{z}^{(k-1)} - \eta^{(k)} \cdot \nabla_{\mathbf{z}} F(\mathbf{z}^{(k-1)})$ ; end return  $\mathbf{z}^{(K)}$ ;

#### Algorithm 1: GRADIENTDESCENT

# Gradient Descent: Convergence

- Letting  $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} F(\mathbf{z})$  denote the global minimum
- Let  $\mathbf{z}^{(k)}$  be our parameter after k updates.
- ▶ Thm: Suppose F is convex and "L-smooth". Using a fixed step size  $\eta \leq \frac{1}{L}$ , we have:

$$F(\mathbf{z}^{(k)}) - F(\mathbf{z}^*) \le \frac{\|\mathbf{z}^{(0)} - \mathbf{z}^*\|^2}{\eta \cdot k}$$

That is the **convergence rate** is  $O(\frac{1}{k})$ .

# Smoothness and Gradient Descent Convergence

 $\blacktriangleright$  Smooth functions: for all z,z'

$$\|\nabla F(z) - \nabla F(z')\| \le L \|z - z'\|$$

- Proof idea:
  - 1. If our gradient is large, we will make good progress decreasing our function value:
  - 2. If our gradient is small, we must have value near the optimal value: