

# Machine Learning (CSE 446): Probabilistic View of Logistic Regression and Linear Regression

Sham M Kakade

© 2018

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

# Announcements

- ▶ Midterm: Weds, Feb 7th. Policies:
  - ▶ You may use a single side of a single sheet of handwritten notes that you prepared.
  - ▶ You must turn your sheet of notes in, with your name on it, in at the conclusion of the exam, even if you never looked at it.
  - ▶ You may not use electronics devices of any sort.
- ▶ Today:  
Review: Regularization and Optimization  
New: (wrap up GD) + probabilistic modeling!

# Review

## Regularization / Ridge Regression

- ▶ **Regularize** the optimization problem:

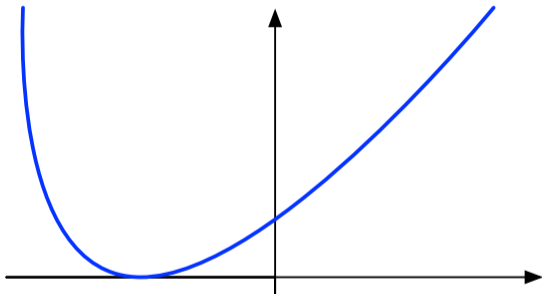
$$\min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w} \cdot \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|^2 = \min_{\mathbf{w}} \frac{1}{N} \|Y - X^T \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

- ▶ This particular case: “Ridge” Regression, Tikhonov regularization
- ▶ The solution is the **least squares estimator**:

$$\mathbf{w}^{\text{least squares}} = \left( \frac{1}{N} X^T X + \lambda \mathbb{I} \right)^{-1} \left( \frac{1}{N} X^T Y \right)$$

**Regularization is often necessary** for the “exact” solution method (regardless of if  $d$  bigger/less than  $N$ )

# Gradient Descent



- ▶ Want to solve:

$$\min_z F(z)$$

- ▶ How should we update  $z$ ?

# Gradient Descent

**Data:** function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$ , number of iterations  $K$ , step sizes  $\langle \eta^{(1)}, \dots, \eta^{(K)} \rangle$

**Result:**  $\mathbf{z} \in \mathbb{R}^d$

initialize:  $\mathbf{z}^{(0)} = \mathbf{0}$ ;

**for**  $k \in \{1, \dots, K\}$  **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

Today

# Gradient Descent: Convergence

- ▶ Denote:  
 $\mathbf{z}^* = \operatorname{argmin}_{\mathbf{z}} F(\mathbf{z})$ : the global minimum  
 $\mathbf{z}^{(k)}$ : 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:

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

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

- ▶ **This Thm applies to both the square loss and logistic loss!**



## Proof intuition: smoothness and GD Convergence

- ▶  $L$ -Smooth functions: “The gradients don't change quickly.” Precisely, For all  $z, z'$

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

- ▶ 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:

## A better idea?

- ▶ Remember the Bayes optimal classifier.  $\mathcal{D}(x, y)$  is the true probability of  $(x, y)$ .

$$\begin{aligned} f^{(\text{BO})}(x) &= \operatorname{argmax}_y \mathcal{D}(x, y) \\ &= \operatorname{argmax}_y \mathcal{D}(y | x) \end{aligned}$$

- ▶ Of course, we don't have  $\mathcal{D}(y | x)$ .

Probabilistic machine learning: **define a probabilistic model** relating random variables  $x$  to  $y$  and **estimate its parameters**.

# A Probabilistic Model for Binary Classification: Logistic Regression

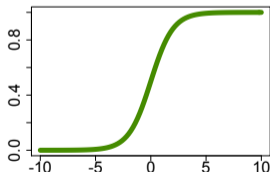
► For  $Y \in \{-1, 1\}$  define  $p_{\mathbf{w},b}(Y | X)$  as:

1. Transform feature vector  $\mathbf{x}$  via the “activation” function:

$$a = \mathbf{w} \cdot \mathbf{x} + b$$

2. Transform  $a$  into a binomial probability by passing it through the logistic function:

$$p_{\mathbf{w},b}(Y = +1 | \mathbf{x}) = \frac{1}{1 + \exp -a} = \frac{1}{1 + \exp -(\mathbf{w} \cdot \mathbf{x} + b)}$$



► If we learn  $p_{\mathbf{w},b}(Y | \mathbf{x})$ , we can (almost) do whatever we like!

# Maximum Likelihood Estimation

The principle of maximum likelihood estimation is to choose our parameters to make our observed data as likely as possible (under our model).

- ▶ Mathematically: find  $\hat{\mathbf{w}}$  that maximizes the probability of the labels  $y_1, \dots, y_n$  given the inputs  $x_1, \dots, x_n$ .
- ▶ Note, by the i.i.d. assumption:

$$\mathcal{D}(y_1, \dots, y_n \mid \mathbf{x}_1, \dots, \mathbf{x}_N) =$$

- ▶ The Maximum Likelihood Estimator (the '**MLE**') is:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^N p_{\mathbf{w}}(y_n \mid \mathbf{x}_n)$$

# Maximum Likelihood Estimation and the Log loss

- ▶ The 'MLE' is:

$$\begin{aligned}\hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w}} \prod_{n=1}^N p_{\mathbf{w}}(y_n | \mathbf{x}_n) \\ &= \operatorname{argmax}_{\mathbf{w}} \log \prod_{n=1}^N p_{\mathbf{w}}(y_n | \mathbf{x}_n) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{n=1}^N \log p_{\mathbf{w}}(y_n | \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N -\log p_{\mathbf{w}}(y_n | \mathbf{x}_n)\end{aligned}$$

- ▶ This is referred to as the **log loss**.

# The MLE for Logistic Regression

- ▶ the MLE for the logistic regression model:

$$\operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N -\log p_{\mathbf{w}}(y_n | \mathbf{x}_n) = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N \log(1 + \exp(-y_n \mathbf{w} \cdot \mathbf{x}_n))$$

- ▶ This is the logistic loss function that we saw earlier.
- ▶ How do we find the MLE?

## Derivation for Log loss for Logistic Regression: scratch space

# Linear Regression as a Probabilistic Model

Linear regression defines  $p_{\mathbf{w}}(Y | X)$  as follows:

1. Observe the feature vector  $\mathbf{x}$ ; transform it via the activation function:

$$\mu = \mathbf{w} \cdot \mathbf{x}$$

2. Let  $\mu$  be the mean of a normal distribution and define the density:

$$p_{\mathbf{w}}(Y | \mathbf{x}) = \frac{1}{\sigma\sqrt{2\pi}} \exp -\frac{(Y - \mu)^2}{2\sigma^2}$$

3. Sample  $Y$  from  $p_{\mathbf{w}}(Y | \mathbf{x})$ .



# Linear Regression-MLE is (Unregularized) Squared Loss Minimization!

$$\operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N -\log p_{\mathbf{w}}(y_n \mid \mathbf{x}_n) \equiv \operatorname{argmin}_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^N \underbrace{(y_n - \mathbf{w} \cdot \mathbf{x}_n)^2}_{\text{SquaredLoss}_n(\mathbf{w}, b)}$$

Where did the variance go?