# CIS 5200: MACHINE LEARNING LINEAR AND LOGISTIC REGRESSION

Content here draws from material by Vatsal Sharan (USC), Christopher De Sa and Kilian Weinberger (Cornell)



19 January 2023

# Surbhi Goel

#### Spring 2023



# LOGISTICS - UPCOMING

#### Homework:

\* HW0 due on Friday, Jan 20, 2023 end of day \* For those on waitlist, email your HWO to Keshav and Wendi (head TAs) \* HWI will be out on Monday, Jan 23, 2023

#### **Recitation:**

Sign up link will be posted on Ed this Friday \* Math background recitation next week **Instructor OH:** 

\* Eric and I will run joint office hours after class on Tuesdays 3:30-4:30

### OUTLINE - TODAY

\* Quick Review of Perceptron \* Logistic Regression \* MLE perspective \* Linear Regression \* Least squares solution \* MLE perspective \* Regularization

### PERCEPTRON - SUMMARY

Input space:  $\mathscr{X} \subseteq \mathbb{R}^d$ **Output space:**  $\mathcal{Y} = \{-1, 1\}$ **Hypothesis Class:**  $\mathcal{F} := \{x \mapsto \operatorname{sign}(w^{\top}x + b) | w \in \mathbb{R}^d, b \in \mathbb{R}\}$ **Loss function:**  $\ell(f(x), y) = \begin{cases} 0 & \text{if } f(x) = y \\ 1 & \text{otherwise.} \end{cases}$ 

**Assumption:** Linearly separable data

**Guarantee:** Zero-error on training data after  $1/\gamma^2$  iterations for margin  $\gamma$ 

# PERCEPTRON - FAILURES

#### Led to the Al winter till mid 1980s **XOR:**

Minsky and Papert in a 1969 book "Perceptrons" showed that Perceptron fails on XOR problems

#### Non-linearly separable data:

Separable in a lifted space

#### Noise:

Hard classifier, cannot model inherent noise

Kernels (later in class)



Non-separable Data

# NON-DETERMINISTIC INPUTS

Perceptron assumed deterministic labels

But there may be inherent uncertainty in the label





<sup>©</sup> Machine Learning @ Berkele

#### We can model this uncertainty using some function $\eta(x) = P(y = 1 | x)$

# LOGISTIC FUNCTION

We can model  $\eta(x) = P(y = 1 | x)$  using different functions

$$\frac{\text{sign}_{0/1}(a)}{\text{Step function}} = \begin{cases} 1 & \text{if } a \ge 0, \\ 0 & \text{otherwise.} \end{cases}$$

$$sigmoid(a) = \frac{1}{1 + \exp(-a)}$$
Sigmoid function

 $P(y = 1 | x) = \eta(x) = \text{sigmoid}(w^{\top}x) = \frac{1}{1 + \exp(-w^{\top}x)}$ 

 $P(y = -1 | x) = 1 - \eta(x) = 1 - \text{sigmoid}(w^{\top}x) = \frac{1}{1 + \exp(w^{\top}x)}$ 



7



### DECISION BOUNDARY

#### How do we decide the label given the logistic model?

$$\frac{P(y = +1 | x)}{P(y = -1 | x)} = \frac{1 + \exp(w^{\top} x)}{1 + \exp(-w^{\top} x)} = e^{-\frac{1}{2}}$$

#### Linear decision boundary







## LOSS FUNCTION

#### Logistic Loss

# $\ell(f(x), y) = \begin{cases} -\log(f(x)) & \text{if } y = 1\\ -\log(1 - f(x)) & \text{otherwise} \end{cases}$

For our setting logistic loss is  $log(1 + exp(-y w^{T}x))$ 

#### 0/I Loss

 $\ell_{0/1}(f(x), y) = \mathbf{1}[f(x) \neq y]$ 

For linear classifier this is  $1[sgn(w^Tx) \neq y] = 1[y w^Tx < 0]$ 

#### Why this loss?



bound of 0/1 loss

9

# PROBABILISTIC VIEW - MAXIMUM LIKELIHOOD ESTIMATOR

Another way to view the supervised learning task is to maximize the likelihood of seeing the training data \* Make an explicit modeling condition on the data distribution \* Find parameters that maximize the probability of seeing the data

> Suppose the parameters of the model are denoted by  $\theta$  $\hat{\mathscr{L}}(\theta) = P(S \mid \theta)$

- S is the training data
- $= P(x_i, y_i \mid \theta)$ Training data is i.i.d.

i=1



# MAXIMUM (CONDITIONAL) LOG LIKELIHOOD Suppose we don't have any assumption on the generation process of x, then we can maximize a conditional likelihood

$$\hat{\mathscr{L}}(\theta) = \prod_{i=1}^{m} P(y_i \mid x_i, \theta)$$

The log-likelihood is then equivalent to:

$$\log \hat{\mathscr{L}}(\theta) = \log \left( \prod_{i=1}^{m} P(y_i \mid x_i, \theta) \right)$$
$$= \sum_{i=1}^{m} \log \left( P(y_i \mid x_i, \theta) \right)$$



log is an increasing function Maximizers of both are identical













# M(C)LE - LOGISTIC REGRESSION

We have the model for P(y | x, w), substituting it gives us



$$\left(P(y_i \mid x_i, w)\right)$$

$$\left(\frac{1}{1 + \exp(-y_i w^{\mathsf{T}} x_i)}\right)$$

$$= -\sum \log \left(1 + \exp(-y_i w^{\mathsf{T}} x_i)\right)$$

#### This is the negative of the logistic loss!

$$\max_{w} \log \hat{\mathscr{L}}(w) = \min_{w} \hat{R}(w)$$

*i*=1

# OGISTIC REGRESSION - TRAINING

### **Training Dataset:** $\mathcal{S} = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\},\$ $x_i \in \mathbb{R}^d, y_i \in \{-1, 1\}$

# **Empirical Risk Minimization:** Find $\hat{w}$ that minimizes $\widehat{R}(w) = \frac{1}{m} \sum_{i=1}^{m} \log_{i=1}^{m} \log_{i=1}^{m$

The problem is convex so we can use convex optimization (will discuss in later lectures)

$$\log\left(1 + \exp(-y_i \ w^{\mathsf{T}} x_i)\right)$$

#### How do we solve this minimization problem?

![](_page_12_Figure_9.jpeg)

# OGISTIC REGRESSION - SUMMARY

Input space:  $\mathscr{X} \subseteq \mathbb{R}^d$  Perceptron **Output space:**  $\mathcal{Y} = [0,1]$   $\mathcal{Y} = \{-1,1\}$ **Hypothesis Class:**  $\mathcal{F} := \{x \mapsto sigmoid(w^T x + b) | w \in \mathbb{R}^d, b \in \mathbb{R}\}$ 

Loss function:  $\ell(f(x), y) = \begin{cases} -\log(f(x)) & \text{if } y = 1\\ -\log(1 - f(x)) & \text{otherwise} \end{cases}$ 

 $\ell(f(x), y) = \begin{cases} 0 & \text{if } f(x) = y \\ 1 & \text{otherwise.} \end{cases}$ 

- $\mathcal{F} := \{ x \mapsto \operatorname{sign}(w^{\mathsf{T}}x + b) \mid w \in \mathbb{R}^d, b \in \mathbb{R} \}$

![](_page_14_Picture_0.jpeg)

#### Predict future outcomes based on past outcomes

Labels  $y \in \mathcal{Y}$ 

 $(\mathcal{Y} = \text{Breeds})$ "Pug" "Chihuahua"

![](_page_14_Figure_5.jpeg)

#### Classification

**Discrete** labels

 $(\mathcal{Y} = \text{Stock prices})$ "\$130.02"

![](_page_14_Picture_9.jpeg)

#### Regression **Continuous** labels

**Task:** Learn predictor  $f: \mathcal{X} \to \mathcal{Y}$ 

![](_page_14_Picture_13.jpeg)

![](_page_14_Picture_14.jpeg)

## HYPOTHESIS CLASS - LINEAR REGRESSORS

Linear regressors  $\mathcal{F} := \{ x \mapsto w^{\mathsf{T}} x + b \mid w \in \mathbb{R}^d, b \in \mathbb{R} \}$ 

![](_page_15_Picture_2.jpeg)

Similar to perceptron, can ignore bias  $\Xi \mathbb{R}^d, b \in \mathbb{R}$ 

![](_page_15_Figure_4.jpeg)

Data from <a href="https://nsidc.org/arcticseaicenews/sea-ice-tools/">https://nsidc.org/arcticseaicenews/sea-ice-tools/</a>

## LOSS FUNCTION

A

![](_page_16_Figure_1.jpeg)

#### **Absolute Loss**

$$\ell(f(x), y) = |f(x) - y|$$
  
bosolute-loss = 
$$\frac{|d_1| + |d_2| + |d_3| + |d_4| + |d_5|}{5}$$

![](_page_16_Figure_5.jpeg)

 $d_{5}$ 

#### How does square loss behave on outliers?

![](_page_16_Picture_9.jpeg)

# LINEAR REGRESSION - TRAINING

- **Empirical Risk Minimization:** Find  $\hat{w}$  that minimizes  $\widehat{R}(w) = \frac{1}{m} \sum_{i=1}^{m} (y_i - w^{\mathsf{T}} x_i)^2$ 
  - How do we solve this minimization problem?

The problem is convex, in fact we can get a closed form solution

**Training Dataset:**  $S = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$ 

#### LEAST SQUARES

Loss is convex  $\Longrightarrow$  differentiate to find minimizer

$$\widehat{R}(w) = \frac{1}{m} \sum_{i=1}^{m} (y_i - w^{\mathsf{T}} x_i)^2$$

![](_page_18_Figure_3.jpeg)

![](_page_18_Figure_5.jpeg)

Least Squares Regression

## SOLVING THE SYSTEM

Normal Equations for Least Squares Regression

If  $X^{\mathsf{T}}X$  is invertible, then

What is the computational cost of computing this?

$$X = \begin{bmatrix} -x_1^{\top} & -\\ -x_2^{\top} & -\\ \vdots & \\ -x_m^{\top} & - \end{bmatrix} \in \mathbb{R}^{m \times d}, Y = \begin{bmatrix} y \\ y \\ y \end{bmatrix}$$

### $X^{\mathsf{T}} X w = X^{\mathsf{T}} Y$

# $\hat{w} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}Y$

 $\hat{Y} = X\hat{w}$  is the projection of Y onto the subspace spanned by the

columns of X,  $\tilde{x}_1, \dots \tilde{x}_d$  where  $\tilde{x}_i = [x_{1i}, \dots, x_{mi}]^\top$ 

Recall that  $X(X^{T}X)^{-1}X$  is the projection matrix on to this subspace

![](_page_19_Picture_15.jpeg)

![](_page_19_Picture_16.jpeg)

# LINEAR REGRESSION - REGULARIZATION

What if  $X^{\dagger}X$  is very close to being singular? This can lead to large values for  $\hat{w}$  which might overfit

# $\widehat{G}(w) = \widehat{R}(w) + \lambda \psi(w)$

 $\psi(w)$  is chosen to be some function that penalizes complexity of w

Common examples include

$$= \frac{1}{m} \sum_{i=1}^{m} (y_i - w^{\mathsf{T}} x_i)^2 + \lambda \psi(w)$$

$$\psi(w) = \|w\|_2^2 \text{ or } \psi(w) = \|w\|_1^2$$

#### RIDGE REGRESSION

$$\widehat{G}(w) = \frac{1}{m} \sum_{i=1}^{m} (y_i - w^{\mathsf{T}} x_i)^2 + \lambda ||w||_2^2$$

$$\hat{w}_{\lambda} = (X^{\mathsf{T}}X + \lambda mI)^{-1}X^{\mathsf{T}}Y$$

Always invertible, eigenvalues are  $\geq \lambda m$ 

![](_page_21_Figure_4.jpeg)

Take derivative and set to 0

![](_page_21_Figure_6.jpeg)

Matrix notation

 $(X^{\top}X + \lambda mI)w = X^{\top}Y$ 

![](_page_21_Picture_10.jpeg)

#### LASSO REGRESSION

![](_page_22_Figure_1.jpeg)

#### Leads to sparsity in the weights!

# Can model as a quadratic program

# LINEAR REGRESSION - SUMMARY

Input space:  $\mathcal{X} \subseteq \mathbb{R}^d$ Output space:  $\mathcal{Y} = \mathbb{R}$ Hypothesis Class:  $\mathcal{F} := \{x \mapsto w^{\mathsf{T}}x + b \mid w \in \mathbb{R}^d, b \in \mathbb{R}\}$ Loss function:  $\ell(f(x), y) = (f(x) - y)^2$ Least Squares solution:  $\hat{w} = (X^{\top}X)^{-1}X^{\top}Y$ 

Next class: Eric will talk about k-nearest neighbors

![](_page_23_Picture_3.jpeg)