Spring 2023!

Introduction to Machine Learning



Lec-Rec 4: Classification & Logistic Regression

Given data
$$D = \{x^{(i)}, y^{(i)}\}_{i=1}^{n}, x^{(i)} \in \mathbb{R}^{d}, y^{(i)} \in \mathbb{R}$$

Define hypothesis class $h(x; \theta) = \theta^{\top} x + \theta_0$

Define cost function
$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} h(x^{(i)}, \theta) - y^{(i)})^2$$

Find best hypothesis

$$\underset{\theta}{\arg\min} J(\theta)$$

Given data
$$D = \{x^{(i)}, y^{(i)}\}_{i=1}^n, x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}$$

Define hypothesis class $h(x; \theta) = \theta^\top x + \theta_0$
 $y^{(i)} \in \{-1, 1\}$

Define cost function

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} h(x^{(i)}, \theta) - y^{(i)})^2$$

Find best hypothesis

$$\underset{\theta}{\arg\min J(\theta)} \qquad ???$$

Example

Regression

- Datum *i*: feature vector $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$ • Label $y^{(i)} \in \mathbb{R}$
- Hypothesis $h: \mathbb{R}^d \to \mathbb{R}$



Example

Regression

- Datum *i*: feature vector
- $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$
 - Label $y^{(i)} \in \mathbb{R}$
- Hypothesis $h: \mathbb{R}^d \to \mathbb{R}$



(Two-class) Classification

- Datum *i*: feature vector $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$
 - Label $y^{(i)} \in \{-1,+1\}$
- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$



Example

Regression

- Datum *i*: feature vector
- $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$
 - Label $y^{(i)} \in \mathbb{R}$
- Hypothesis $h: \mathbb{R}^d \to \mathbb{R}$



(Two-class) Classification

- Datum *i*: feature vector $x^{(i)} = (x_1^{(i)}, \dots, x_d^{(i)})^\top \in \mathbb{R}^d$
 - Label $y^{(i)} \in \{-1,+1\}$
- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$



- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$
- Linear classifier: a hypothesis class which labels +1 on one side of a hyperplane and -1 on the other

- Hypothesis $h : \mathbb{R}^d \to \{-1, +1\}$
- Linear classifier: a hypothesis class which labels +1 on one side of a hyperplane and -1 on the other
- Hyperplane: high-dimensional line that can be characterized by:

$$\theta^{\top} x + \theta_0 = 0$$















$$\theta^{\top} x = 0$$

• θ is vector that is **normal** (or orthogonal) to the hyperplane



$$\theta^{\top} x = 0$$

- θ is vector that is **normal** (or orthogonal) to the hyperplane
- $heta_0$ controls offset from origin

 $\theta^{\top} x + \theta_0 = 0$



• Hyperplane:

$$\theta^{\top} x + \theta_0 = 0$$

• Decision rule:

$$\begin{split} \mathfrak{n}(\mathbf{x}; \theta, \theta_0) &= \operatorname{sign}(\theta^{\mathsf{T}} \mathbf{x} + \theta_0) \\ &= \begin{cases} +1 & \text{if } \theta^{\mathsf{T}} \mathbf{x} + \theta_0 > 0 \\ -1 & \text{otherwise} \end{cases} \end{split}$$

y = Wearing a coat?







Direction along x1 axis Direction along x2 axis





Classifier Performance

 x_1

Which classifier is better? y = Wearing a coat? y = Wearing a coat? $x^{(3)}$ $x^{(3)}$ $x^{(2)}$ x_2 $x^{(1)}$ x_2 $x^{(1)}$ eed (kph) Wind speed (kpt Win Temperature (C) x_1 Temperature (C)

Classifier Performance

- Which classifier is better?
- "Zero-one" loss:

$$L(g,a) = \begin{cases} 0 \text{ if } g = a \\ 1 \text{ else} \end{cases} \begin{array}{c} \texttt{g: guess,} \\ \texttt{a: actual} \end{cases}$$

• Dataset level error: 1 - accuracy

$$\frac{1}{n}\sum_{i=1}^{n} L(g^{(i)}, y^{(i)})$$

• Issue 1: No notion of uncertainty



• Issue 1: No notion of uncertainty





• Issue 2: Gradient descent not possible

$$\begin{split} J(\theta) &= \frac{1}{n} \sum_{i=1}^{n} L(g^{(i)}, y^{(i)}) \\ &= \frac{1}{n} \sum_{i=1}^{n} \mathbbm{1} \left\{ y^{(i)} \neq \mathsf{sign}(\theta^{\top} x^{(i)} + \theta_0) \right\} \end{split}$$

What is the gradient of the above with respect the model parameters?

• Issue 2: Gradient descent not possible

$$\begin{split} J(\theta) &= \frac{1}{n} \sum_{i=1}^{n} L(g^{(i)}, y^{(i)}) \\ &= \frac{1}{n} \sum_{i=1}^{n} \mathbbm{1} \left\{ y^{(i)} \neq \mathsf{sign}(\theta^{\top} x^{(i)} + \theta_0) \right\} \end{split}$$

• Gradient is zero "almost everywhere"

$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta} \mathbb{1} \left\{ y^{(i)} \neq \mathsf{sign}(\theta^{\top} x^{(i)} + \theta_0) \right\}$$

























Cost Function

- What (differentiable) cost function should we optimize?
- Mean squared error?

$$J(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left(\sigma(\theta^{\top} x^{(i)} + \theta_0) - y^{(i)} \right)^2$$
$$y \in \{0, 1\}$$

• Will "work", but loss is non-convex and somewhat unnatural.

- Idea: find parameters that would **maximize the probability of observed data.**
- Probability of a label

$$P(Y = +1 | X = x) = \sigma(\theta^{\top} x + \theta_0)$$

"Probability that the label is +1 given datum x"

- Idea: find parameters that would **maximize the probability of observed data.**
- Probability of a label

 $P(Y = +1 | X = x) = \sigma(\theta^{\top} x + \theta_0)$ $P(Y = -1 | X = x) = 1 - \sigma(\theta^{\top} x + \theta_0)$ $P(Y = y | X = x) = \sigma(\theta^{\top} x + \theta_0)^{\mathbb{1}\{y = +1\}} \cdot (1 - \sigma(\theta^{\top} x + \theta_0))^{\mathbb{1}\{y = -1\}}$

- Idea: find parameters that would **maximize the probability of observed data.**
- Probability of a label

 $P(Y = +1 | X = x) = \sigma(\theta^{\top} x + \theta_0)$ $P(Y = -1 | X = x) = 1 - \sigma(\theta^{\top} x + \theta_0)$ $P(Y = y | X = x) = \sigma(\theta^{\top} x + \theta_0)^{\mathbb{1}\{y = +1\}} \cdot (1 - \sigma(\theta^{\top} x + \theta_0))^{\mathbb{1}\{y = -1\}}$

• Probability of all data: $P(Data) = \prod_{i=1}^{n} P(Y = y^{(i)} | X = x^{(i)})$ "Likelihood" function

- Idea: find parameters that would **maximize the probability of observed data.**
- Maximum Likelihood Estimation

$$\arg \max_{\theta} P(\mathsf{Data})$$

$$= \arg \max_{\theta} \log P(\mathsf{Data})$$

$$= \arg \min_{\theta} -\log P(\mathsf{Data})$$

$$= \arg \min_{\theta} -\frac{1}{n} \log P(\mathsf{Data})$$

Average negative log likelihood

Cost function: average negative log likelihood

$$\begin{split} J(\theta) &= -\frac{1}{n} \log P(\mathsf{Data}) = -\frac{1}{n} \log \prod_{i=1}^{n} P(Y = y^{(i)} | X = x^{(i)}) \\ &= -\frac{1}{n} \sum_{i=1}^{n} \log P(Y = y^{(i)} | X = x^{(i)}) \\ &= -\frac{1}{n} \sum_{i=1}^{n} \log \left(\sigma(\theta^{\top} x^{(i)} + \theta_0)^{\mathbb{1} \left\{ y^{(i)} = +1 \right\}} \cdot (1 - \sigma(\theta^{\top} x^{(i)} + \theta_0))^{\mathbb{1} \left\{ y^{(i)} = -1 \right\}} \\ &= -\frac{1}{n} \sum_{i=1}^{n} \mathbb{1} \left\{ y^{(i)} = +1 \right\} \log \sigma(\theta^{\top} x^{(i)} + \theta_0) + \begin{array}{c} \text{Differentiable (and convex)} \\ & \text{with respect to model} \\ \mathbb{1} \left\{ y^{(i)} = -1 \right\} \log(1 - \sigma(\theta^{\top} x^{(i)} + \theta_0)) \end{array}$$

Logistic Regression

- Logistic Regression
 - Probability of label given by sigmoid function
 - Trained to maximize likelihood

Uncertainty-awareness + differentiable loss



Cost function: average negative log likelihood

$$y \in \{-1, +1\} \qquad J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} \mathbb{1} \left\{ y^{(i)} = +1 \right\} \log \sigma(\theta^{\top} x^{(i)} + \theta_0) + \mathbb{1} \left\{ y^{(i)} = -1 \right\} \log(1 - \sigma(\theta^{\top} x^{(i)} + \theta_0))$$

$$\begin{split} y \in \{0, 1\} & J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log \sigma(\theta^{\top} x^{(i)} + \theta_0) + \\ \text{Use this formulation} & (1 - y^{(i)}) \log(1 - \sigma(\theta^{\top} x^{(i)} + \theta_0)) \\ \text{from now on} \end{split}$$

Logistic Regression Gradients

$$y^{(i)} \in \{0, 1\} \qquad g^{(i)} = \frac{1}{1 + \exp(-(\theta^{\top} x^{(i)} + \theta_0))}$$
$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log(1 - g^{(i)})$$
$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - g^{(i)}) x^{(i)}$$

Try to derive above and also find $\nabla_{\theta_0} J(\theta)$

Logistic vs. Linear Regression Gradients

$$\begin{split} & \text{Logistic} \\ & \text{Regression} \\ & \begin{array}{l} J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log(1 - g^{(i)}) \\ & \nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)}) x^{(i)} = -\frac{1}{n} \sum_{i=1}^{n} \left(y^{(i)} - \sigma(\theta^{\top} x^{(i)} + \theta_0) \right) x^{(i)} \\ & \\ & \text{Linear} \\ & \text{Regression} \\ \end{array} \\ & \begin{array}{l} J(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)})^2 \\ & \nabla_{\theta} J(\theta) = -\frac{2}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)}) x^{(i)} = -\frac{2}{n} \sum_{i=1}^{n} \left(y^{(i)} - (\theta^{\top} x^{(i)} + \theta_0) \right) x^{(i)} \\ \end{array} \end{split}$$

Logistic vs. Linear Regression Gradients

$$\begin{split} & \text{Logistic}_{\substack{\mathsf{Regression}\\\mathsf{Regression}}} & J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log(1 - g^{(i)}) \\ & \nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)}) x^{(i)} = -\frac{1}{n} \sum_{i=1}^{n} \left(y^{(i)} - \sigma(\theta^{\top} x^{(i)} + \theta_0) \right) x^{(i)} \\ & \text{Linear}_{\substack{\mathsf{Regression}\\\mathsf{Regression}}} & J(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)})^2 \\ & \nabla_{\theta} J(\theta) = -\frac{2}{n} \sum_{i=1}^{n} (y^{(i)} - g^{(i)}) x^{(i)} = -\frac{2}{n} \sum_{i=1}^{n} \left(y^{(i)} - (\theta^{\top} x^{(i)} + \theta_0) \right) x^{(i)} \end{split}$$

Intuition: Adjust model parameters to point in the same direction as datum, weighted by how "wrong" the current model is (residual)

Classification rule:

$$h(x) = \mathbb{1}\{\sigma(\theta^{\top}x + \theta_0) > 0.5\} = \mathbb{1}\{\theta^{\top}x + \theta_0 > 0\}$$

⇒ Still a linear classifier! Logistic regression is just one way of learning a hyperplane

Evaluation:

- Accuracy on test set
- Negative log likelihood of test set

