

10-301/601: Introduction to Machine Learning

Lecture 10 – Logistic Regression

Henry Chai

6/6/23

Front Matter

- Announcements:
 - None!
- Recommended Readings:
 - Murphy, [Chapters 8.1-8.3](#)

Recall: Probabilistic Learning

- Previously:
 - (Unknown) Target function, $c^*: \mathcal{X} \rightarrow \mathcal{Y}$
 - Classifier, $h: \mathcal{X} \rightarrow \mathcal{Y}$
 - Goal: find a classifier, h , that best approximates c^*
- Now:
 - (Unknown) Target *distribution*, $y \sim P^*(Y|\mathbf{x})$
 - Distribution, $P(Y|\mathbf{x})$
 - Goal: find a distribution, \underbrace{P} , that best approximates P^*

Building a Probabilistic Classifier

- Define a decision rule
 - Given a test data point \mathbf{x}' , predict its label \hat{y} using the *posterior distribution* $P(Y = y|X = \mathbf{x}')$
 - Common choice: $\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y|X = \mathbf{x}')$
- Model the posterior distribution
 - Option 1 - Model $P(Y|X)$ directly as some function of X (today!)
 - Option 2 - Use Bayes' rule (later):

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto P(X|Y) P(Y)$$

Modelling the Posterior

- Suppose we have binary labels $y \in \{0,1\}$ and D -dimensional inputs $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- **Assume**

$$P(Y = 1|\mathbf{x}) = \text{logit}(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$
$$= \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$$

- This implies two useful facts:

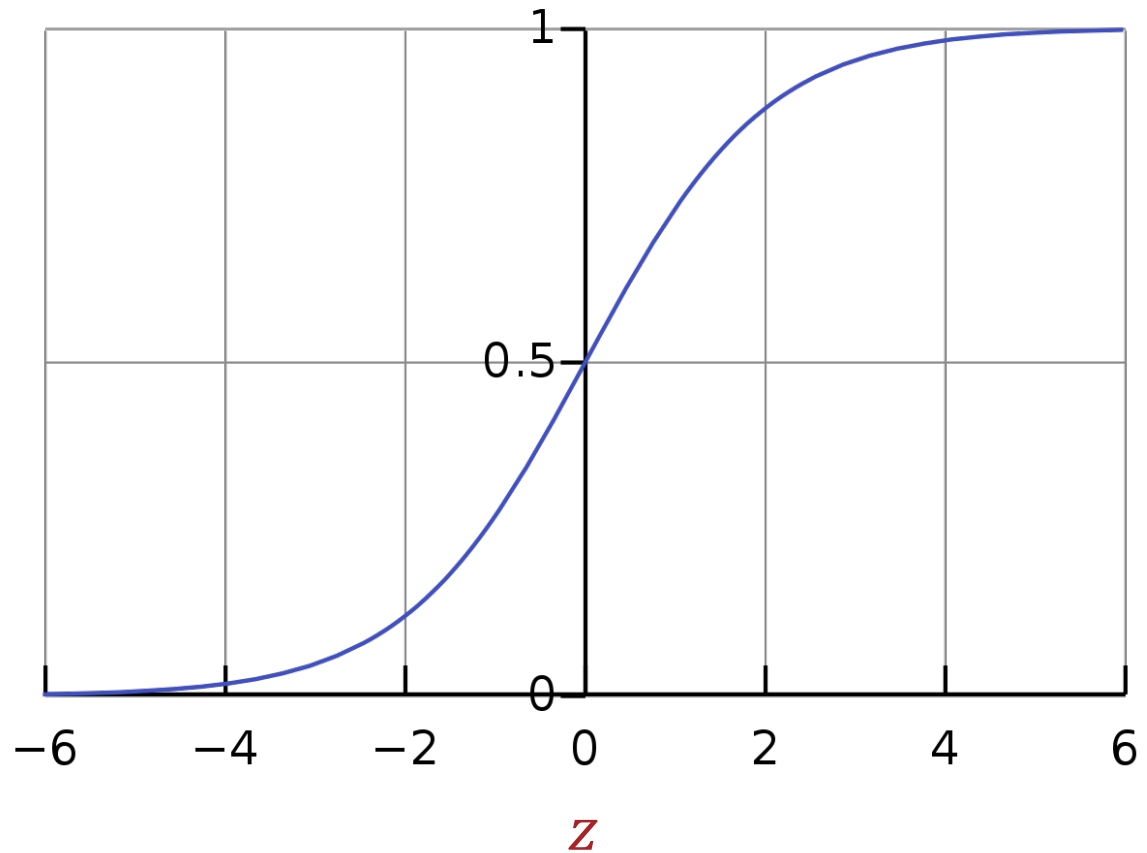
$$1. P(Y=0|\mathbf{x}) = 1 - P(Y=1|\mathbf{x}) = 1 - \frac{\exp(\mathbf{w}^T \mathbf{x})}{\exp(\mathbf{w}^T \mathbf{x}) + 1} = \frac{1}{\exp(\mathbf{w}^T \mathbf{x}) + 1}$$

$$2. \frac{P(Y=1|\mathbf{x})}{P(Y=0|\mathbf{x})} = \exp(\mathbf{w}^T \mathbf{x}) \rightarrow \log \text{ odds} = \log(\exp(\mathbf{w}^T \mathbf{x}))$$

$= \mathbf{w}^T \mathbf{x}$
is linear (in \mathbf{x})

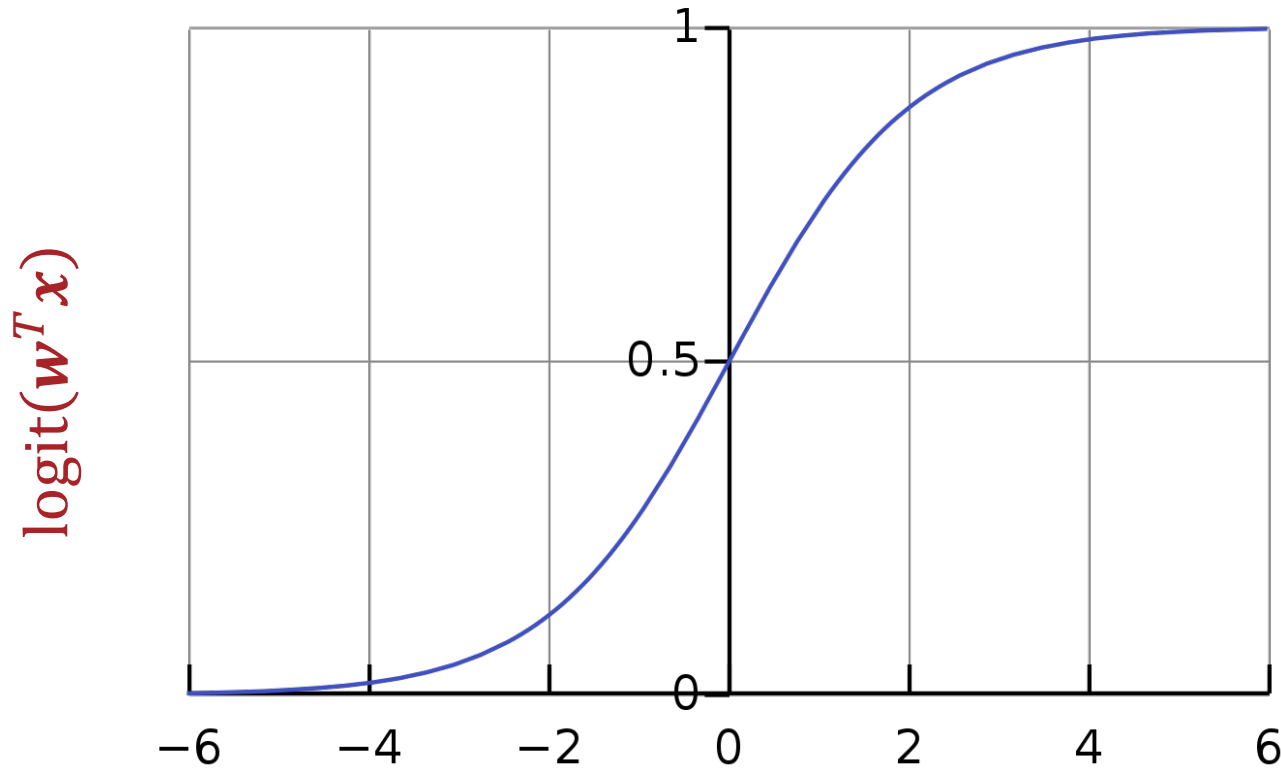
Logistic Function

$$\text{logit}(z) = \frac{1}{1 + e^{-z}}$$



Why use the Logistic Function?

- gives rise to a linear decision boundary



- $\text{logit} : \mathbb{R} \mapsto [0, 1]$ $w^T x$

- differentiable everywhere \Rightarrow convenient optimization properties

- centered at 0.5 $\Rightarrow \text{logit}(0) = 0.5$

Logistic Regression Decision Boundary

$$\hat{y} = \begin{cases} 1 & \text{if } P(Y = 1|\mathbf{x}) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

why $\frac{1}{2}$?

$$P(Y=1|\mathbf{x}) = \frac{1}{1 + \exp(-w^T \mathbf{x})} \geq \frac{1}{2}$$

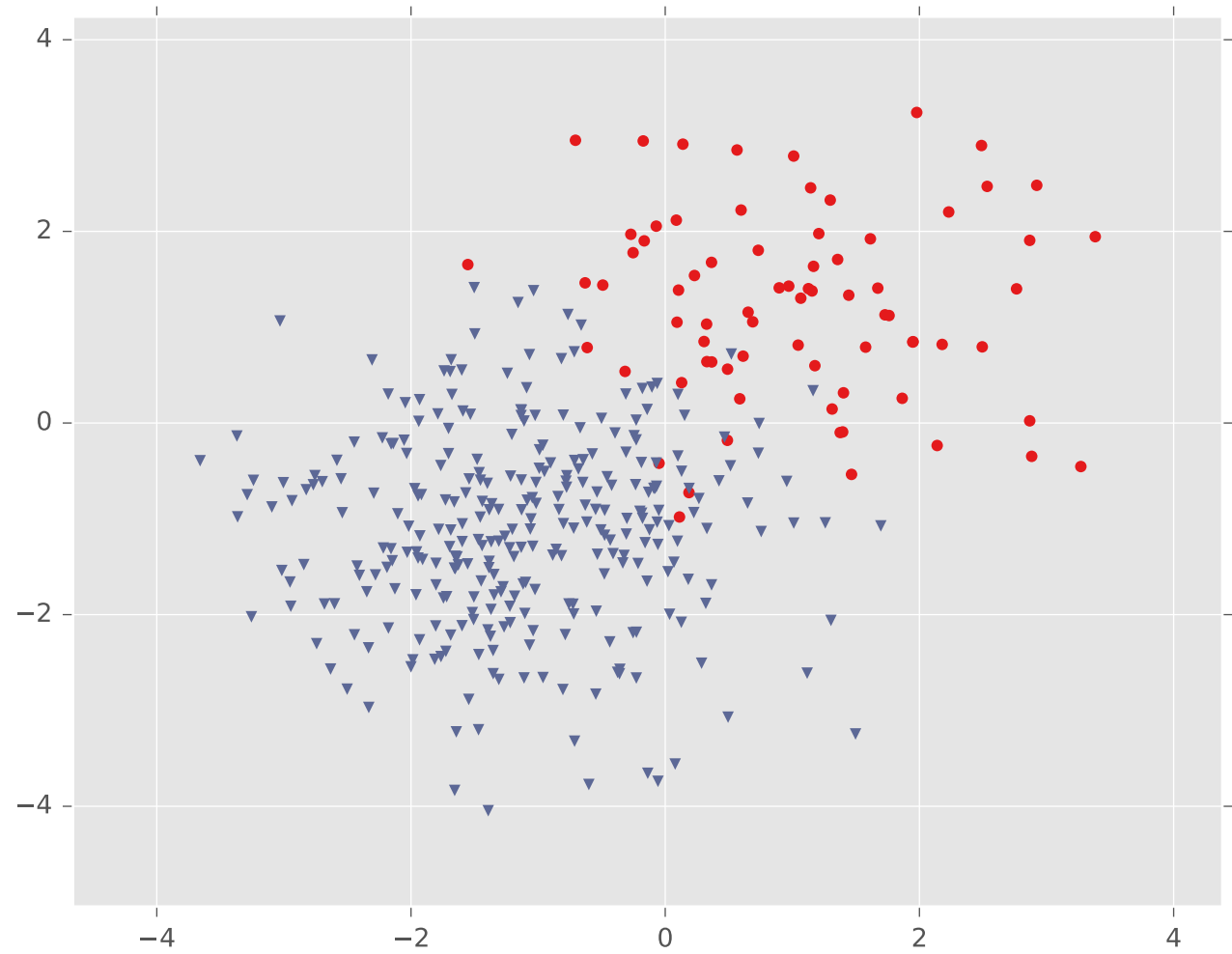
$$\rightarrow 2 \geq 1 + \exp(-w^T \mathbf{x})$$

$$\rightarrow 1 \geq \exp(-w^T \mathbf{x})$$

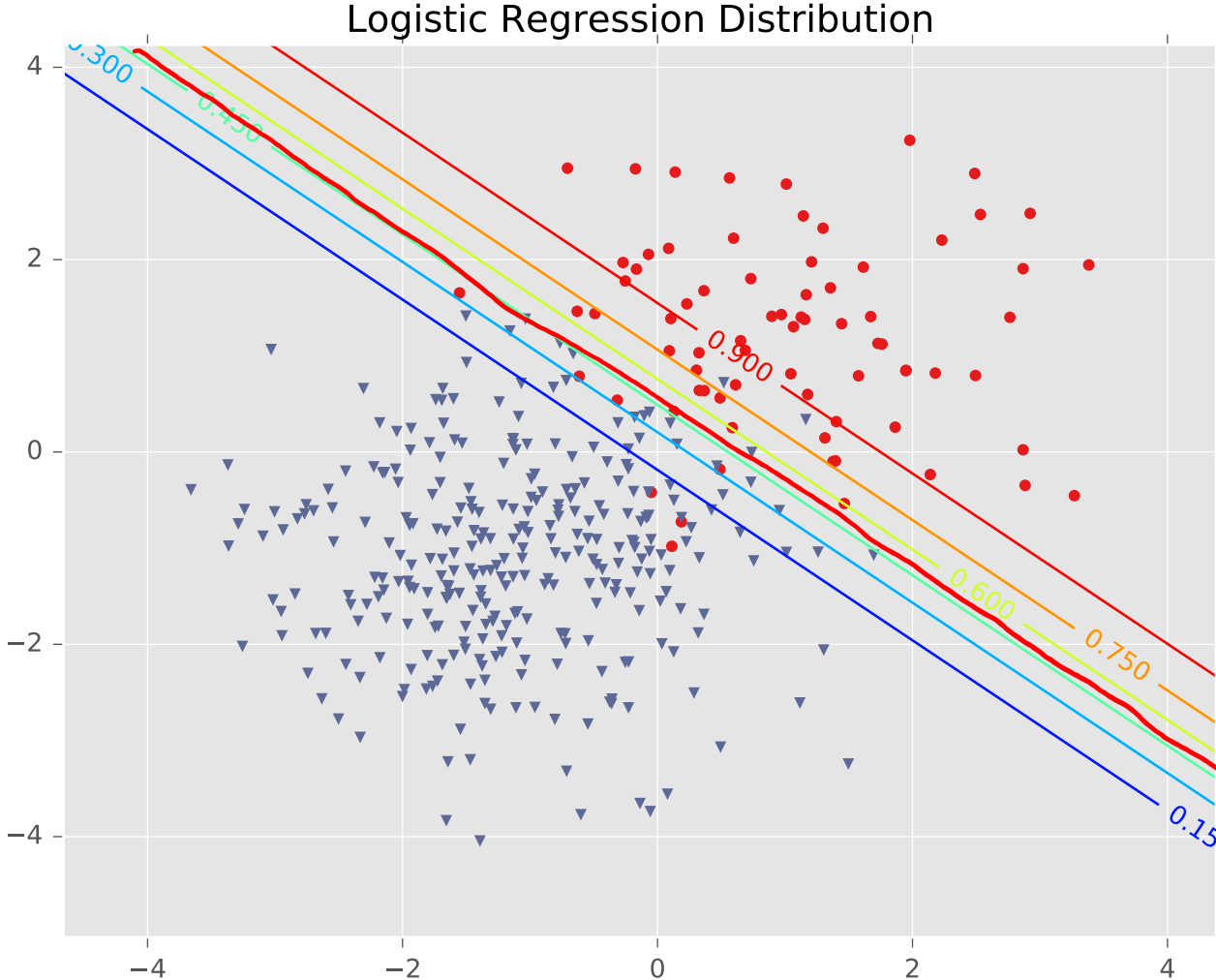
$$\rightarrow \ln(1) = 0 \geq -w^T \mathbf{x}$$

$$\rightarrow w^T \mathbf{x} \geq 0$$

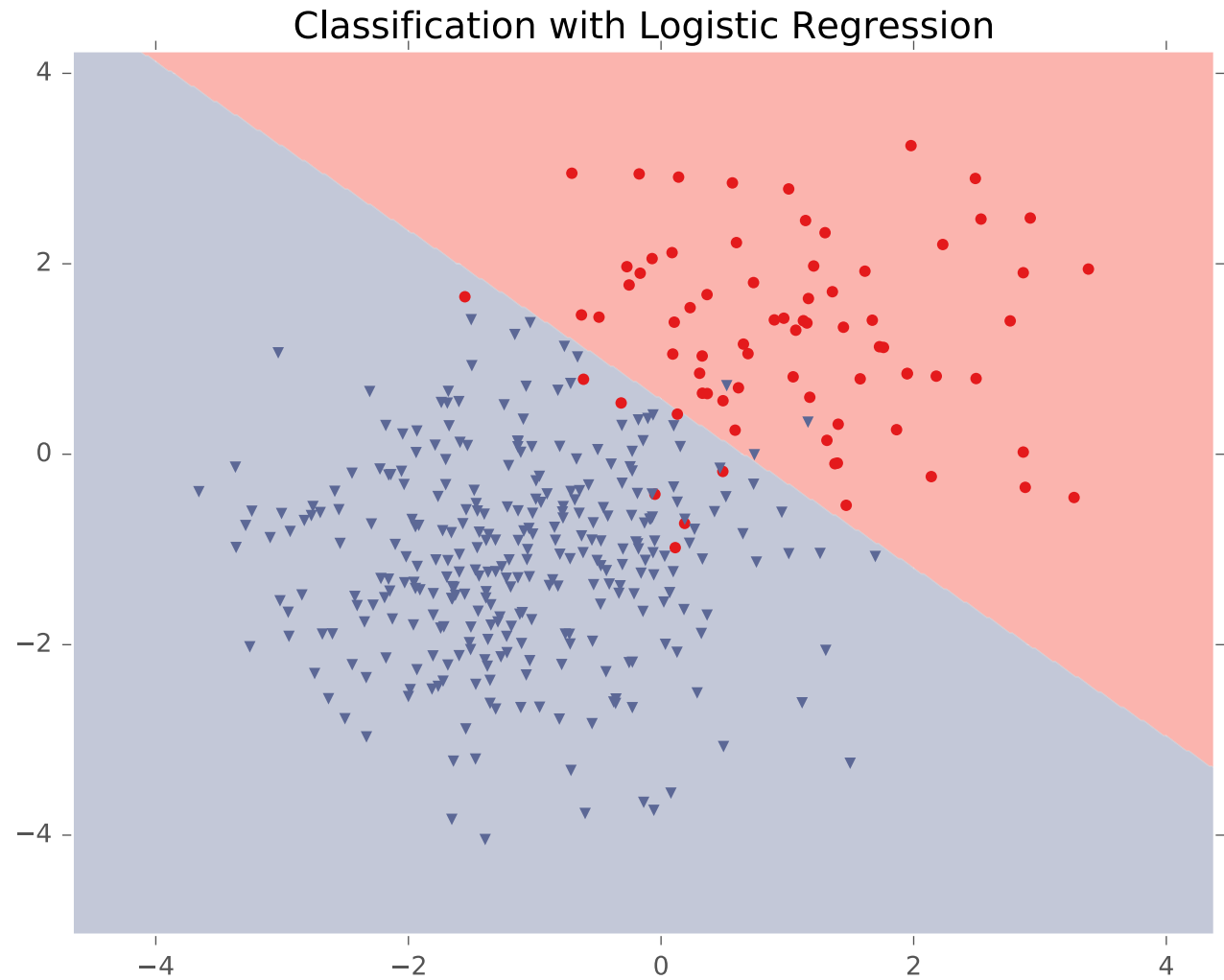
Logistic Regression Decision Boundary



Logistic Regression Decision Boundary



Logistic Regression Decision Boundary



General Recipe for Machine Learning

- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for Logistic Regression

- Define a model and model parameters
 - Assume that the data is i.i.d.
 - Assume $P(Y=1|x) = \text{logit}(w^T x)$
 - Parameters: $w = [w_0, w_1, \dots, w_D]^T$
- Write down an objective function
 - Maximize the log-likelihood
 - Minimize the negative "conditional" log-likelihood
- Optimize the objective w.r.t. the model parameters
 - ???

Full or joint likelihood

$$\rightarrow \ell_D^{\text{full}}(\mathbf{w}) = P(\underbrace{x^{(1)}}_{\mathbf{w}}, y^{(1)}, \underbrace{x^{(2)}}_{\mathbf{w}}, y^{(2)}, \dots, \underbrace{x^{(N)}}_{\mathbf{w}}, y^{(N)})$$

Find \mathbf{w} that minimizes

$$\begin{aligned} \ell_D(\mathbf{w}) &= -\log P(y^{(1)}, \dots, y^{(N)} | x^{(1)}, \dots, x^{(N)}, \mathbf{w}) = -\log \prod_{n=1}^N P(y^{(n)} | x^{(n)}, \mathbf{w}) \\ &= -\ln \prod_{n=1}^N P(Y=1 | x^{(n)}, \mathbf{w})^{y^{(n)}} P(Y=0 | x^{(n)}, \mathbf{w})^{1-y^{(n)}} \\ &= -\sum_{n=1}^N y^{(n)} \ln (P(Y=1 | x^{(n)}, \mathbf{w})) + \underbrace{(1-y^{(n)})}_{\ln(P(Y=0 | x^{(n)}, \mathbf{w}))} \\ &= -\sum_{n=1}^N y^{(n)} \ln \left(\frac{P(Y=1 | x^{(n)}, \mathbf{w})}{P(Y=0 | x^{(n)}, \mathbf{w})} \right) + (1) \ln(P(Y=0 | x^{(n)}, \mathbf{w})) \\ &= -\sum_{n=1}^N y^{(n)} (\mathbf{w}^T x^{(n)}) + \ln \left(\frac{1}{1 + \exp(\mathbf{w}^T x^{(n)})} \right) \\ &= -\sum_{n=1}^N y^{(n)} (\mathbf{w}^T x^{(n)}) - \ln (1 + \exp(\mathbf{w}^T x^{(n)})) \end{aligned}$$

Setting the Parameters via Maximum Negative Log-Likelihood Estimation! (MCLE)

Minimizing the Negative Conditional (log-)Likelihood

$$\hookrightarrow \ell_D(w) = - \sum_{n=1}^N y^{(n)} (w^T x^{(n)}) - \ln(1 + \exp(w^T x^{(n)}))$$

$$\nabla_w \ell_D(w) = - \sum_{n=1}^N \nabla_w \left(y^{(n)} (w^T x^{(n)}) - \ln(1 + \exp(w^T x^{(n)})) \right)$$

$$= - \sum_{n=1}^N \underbrace{y^{(n)}}_{\substack{\text{constant} \\ \text{w.r.t. } w}} \underbrace{x^{(n)}}_{\substack{\text{constant} \\ \text{w.r.t. } w}} - \frac{1}{1 + \exp(w^T x^{(n)})} \exp(w^T x^{(n)}) \underbrace{x^{(n)}}_{\substack{\text{constant} \\ \text{w.r.t. } w}}$$

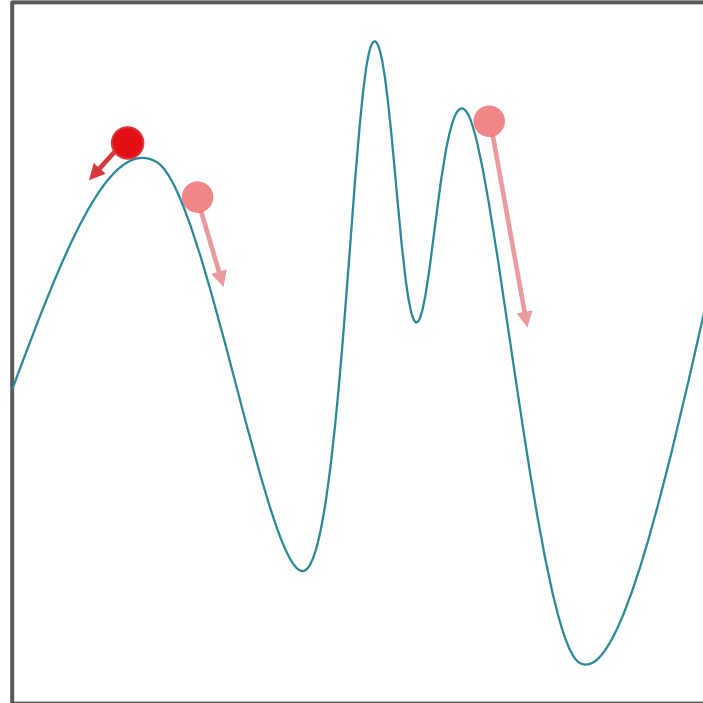
$$= - \underbrace{\sum_{n=1}^N x^{(n)} y^{(n)}}_N - \underbrace{\frac{\exp(w^T x^{(n)})}{1 + \exp(w^T x^{(n)})}}_N$$

$$= - \sum_{n=1}^N x^{(n)} \left(y^{(n)} - P(Y=1 | x^{(n)}, w) \right)$$

Hw $\ell_D(w)$ is positive semi-definite

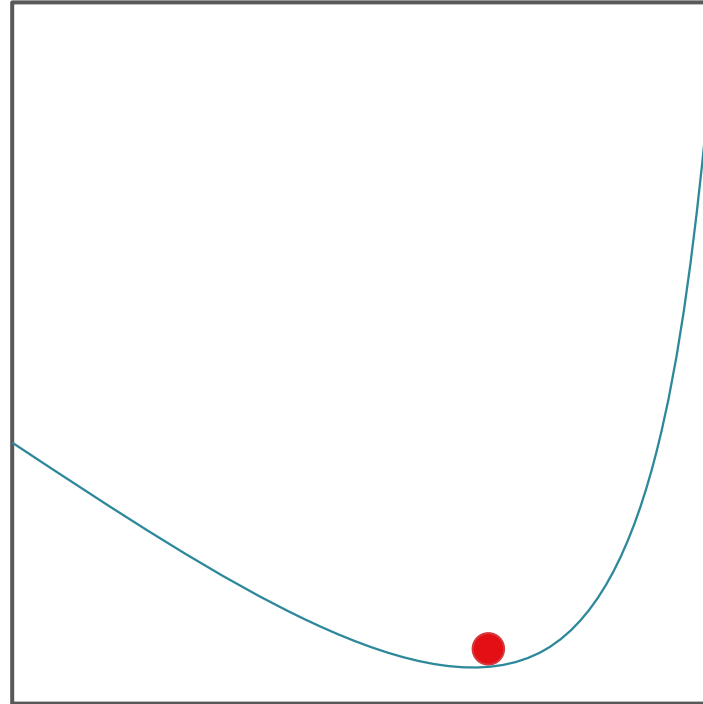
Recall: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



Recall: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



- Good news: the negative conditional log-likelihood, like the squared error, is also convex!

Gradient Descent

• Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta^{(0)}$

1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$

2. While TERMINATION CRITERION is not satisfied

a. Compute the gradient:

$$O(ND) \left\{ \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) = \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \right.$$

b. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)})$

c. Increment t : $t \leftarrow t + 1$

• Output: $\mathbf{w}^{(t)}$

Stochastic Gradient Descent

- Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \underline{\eta}_{SGD}^{(0)}$
- 1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
 - a. Randomly sample a data point from \mathcal{D} , $(\mathbf{x}^{(n)}, y^{(n)})$
 - b. Compute the pointwise gradient:
$$\underline{\nabla}_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) = \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)})$$
 - c. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{SGD}^{(0)} \underline{\nabla}_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)})$
 - d. Increment t : $t \leftarrow t + 1$
- Output: $\mathbf{w}^{(t)}$

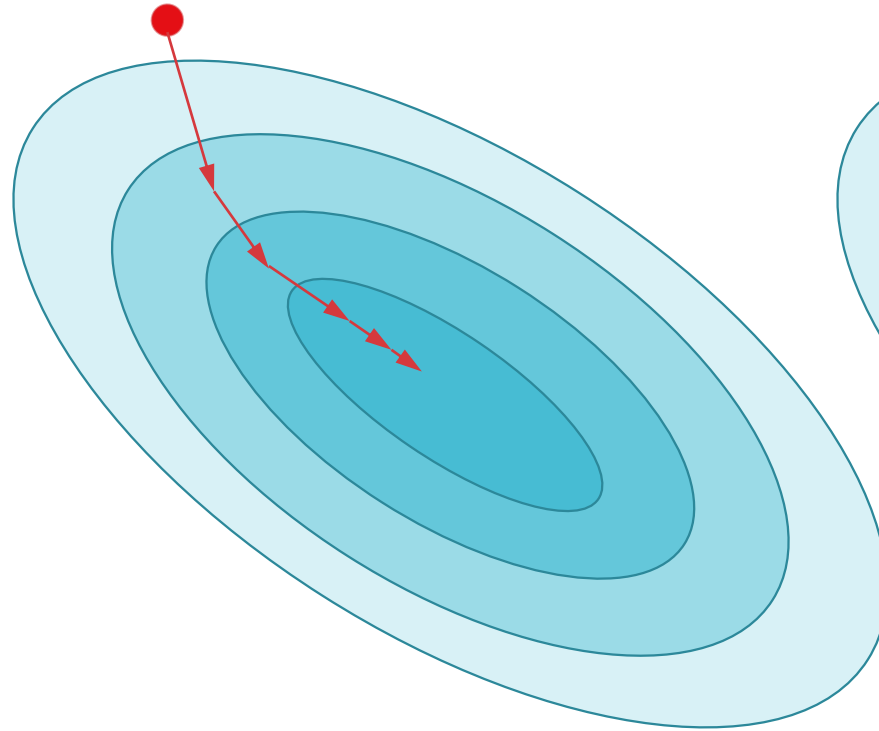
Stochastic Gradient Descent

- If the data point is sampled uniformly at random, then the expected value of the pointwise gradient is proportional to the full gradient:

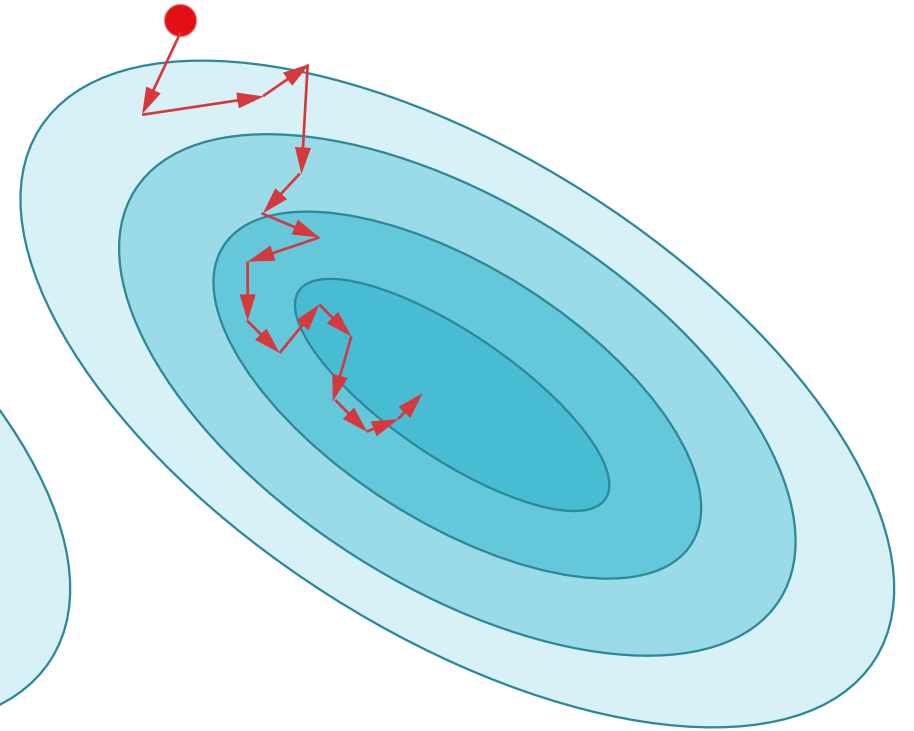
$$\begin{aligned} E \left[\nabla_{\mathbf{w}} \ell_{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}}(\mathbf{w}^{(t)}) \right] &= \frac{1}{N} \sum_{n=1}^N \nabla_{\mathbf{w}} \ell^{(n)}(\mathbf{w}^{(t)}) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (P(Y = 1 | \mathbf{x}^{(n)}, \mathbf{w}^{(t)}) - y^{(n)}) \\ &= \frac{1}{N} \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}^{(t)}) \end{aligned}$$

- In practice, the data set is randomly shuffled then looped through so that each data point is used equally often

Stochastic Gradient Descent vs. Gradient Descent



Gradient Descent



Stochastic Gradient Descent

Mini-batch Stochastic Gradient Descent

• Input: $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N, \eta_{MB}^{(0)}, B$

1. Initialize $\mathbf{w}^{(0)}$ to all zeros and set $t = 0$
2. While TERMINATION CRITERION is not satisfied

a. Randomly sample B data points from \mathcal{D} :

$$\mathcal{D}_{batch} \{(\mathbf{x}^{(b)}, y^{(b)})\}_{b=1}^B$$

b. Compute the gradient w.r.t. the sampled *batch*:

$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)}) = \sum_{b=1}^B \mathbf{x}^{(b)} (P(Y = 1 | \mathbf{x}^{(b)}, \mathbf{w}) - y^{(b)})$$


c. Update \mathbf{w} : $\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta_{MB}^{(0)} \nabla_{\mathbf{w}} \ell_{\mathcal{D}_{batch}}(\mathbf{w}^{(t)})$

d. Increment t : $t \leftarrow t + 1$

• Output: $\mathbf{w}^{(t)}$

Key Takeaways

$$f(x) = x^2 \quad \frac{\partial f}{\partial x} = 2x \quad \frac{\partial^2 f}{\partial x^2} = 2$$



$$H_{\ell} = \begin{bmatrix} \frac{\partial^2 \ell}{\partial v_1^2} & \frac{\partial^2 \ell}{\partial v_1 \partial v_2} & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

- Logistic regression
 - Logistic function induces a linear decision boundary
 - Conditional likelihood maximization
- Gradient descent vs. stochastic gradient descent tradeoffs

$$= \begin{bmatrix} \frac{\partial \ell}{\partial v_1} & \frac{\partial \ell}{\partial v_2} & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$V \vec{s}, \quad \vec{s}^T H \vec{s} \geq 0$$