

CS 360: Machine Learning

Sara Mathieson, Sorelle Friedler

Spring 2024



Haverford
COLLEGE

Admin

- Checkpoint in lab today: finished **Part 1**
- Office hour change: **Wednesday 3-4pm in Zubrow** (Sara)
 - No office hours on Thursday! (faculty meeting)
- **Lab 4** due Thursday Feb 22

Outline for Feb 20

- Stochastic gradient descent
- Maximizing likelihood functions
- Logistic regression likelihood
- Extending logistic regression to multi-class classification (softmax)
- Regularization

Outline for Feb 20

- Stochastic gradient descent
- Maximizing likelihood functions
- Logistic regression likelihood
- Extending logistic regression to multi-class classification (softmax)
- Regularization

High-level idea of gradient descent

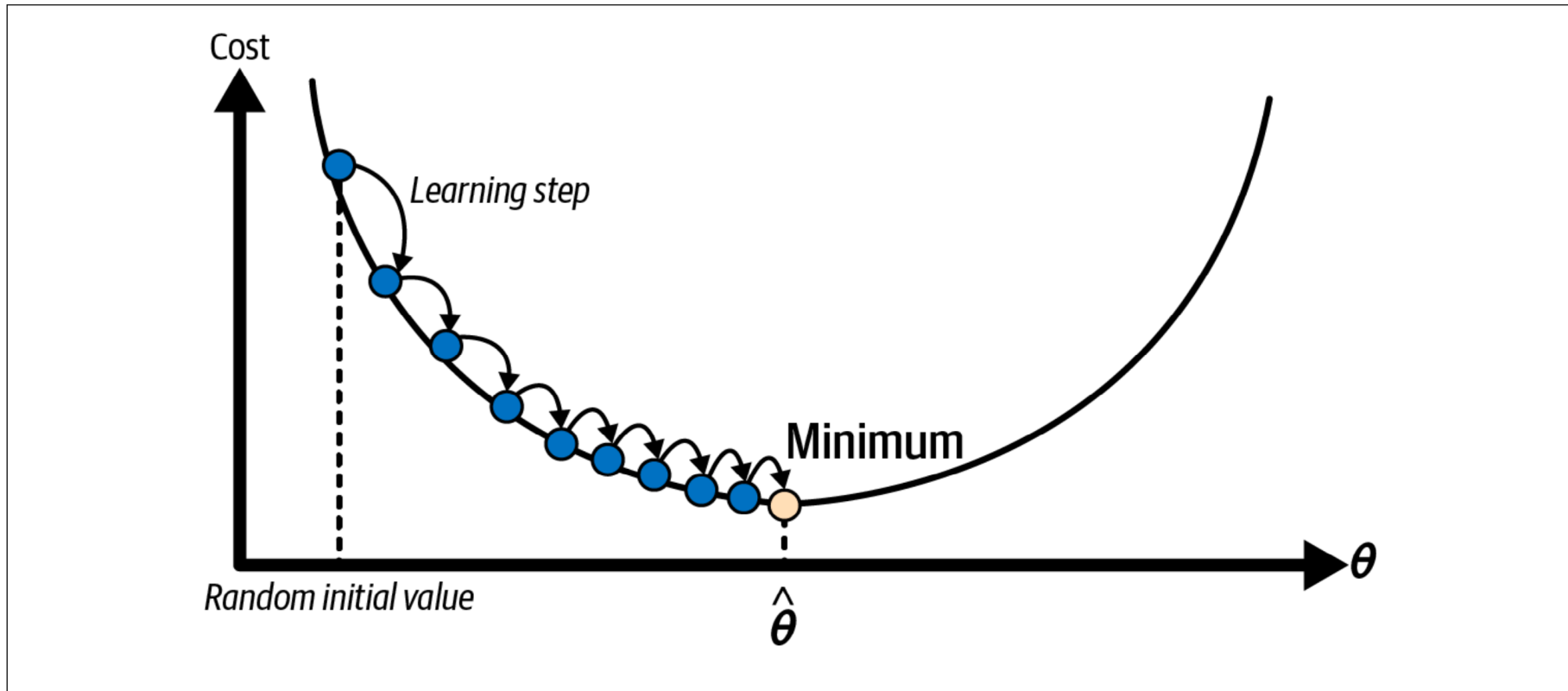
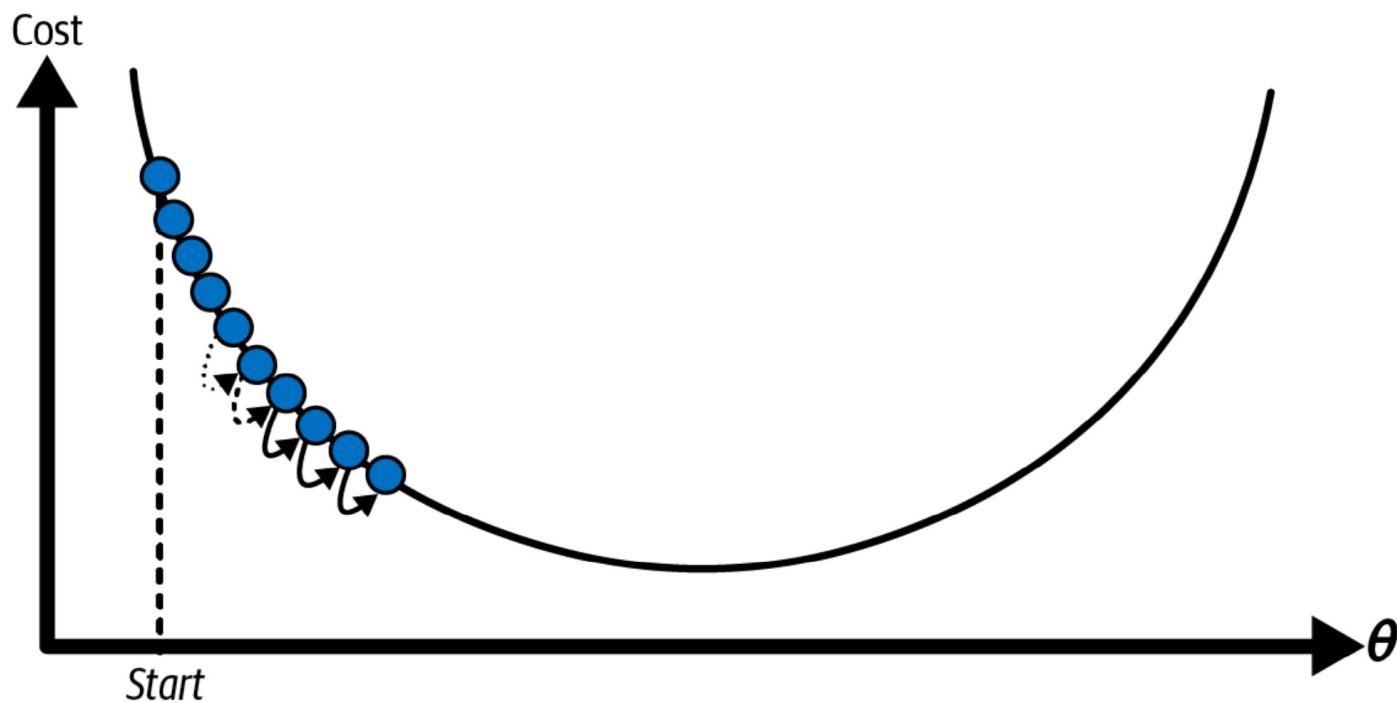
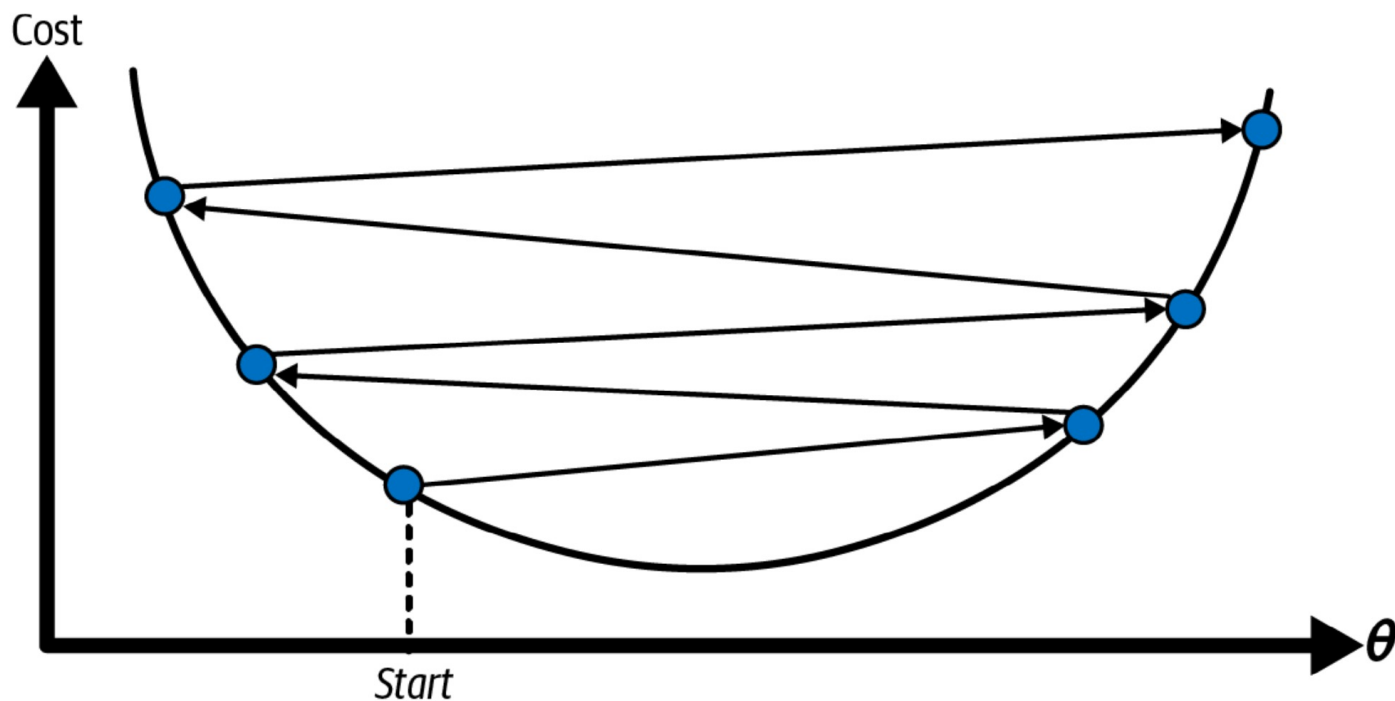


Figure 4-3. In this depiction of gradient descent, the model parameters are initialized randomly and get tweaked repeatedly to minimize the cost function; the learning step size is proportional to the slope of the cost function, so the steps gradually get smaller as the cost approaches the minimum

Gradient descent: learning rate too small



Gradient descent: learning rate too high



Gradient descent: no global optima

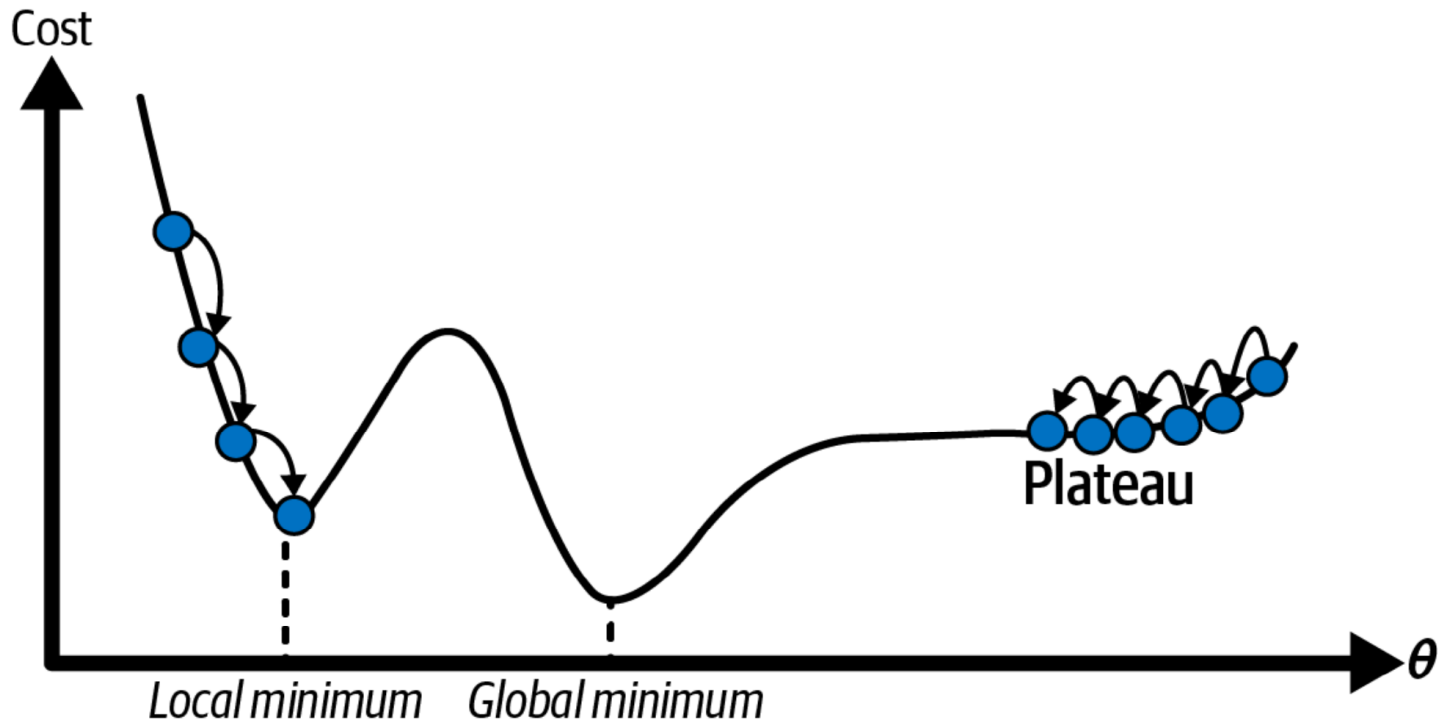


Figure 4-6. Gradient descent pitfalls

Stochastic Gradient Descent (high-level)

set $\mathbf{w} = \mathbf{0}$ vector

while **cost** $J(\mathbf{w})$ still changing (or max iter reached):

 shuffle data points

 for $i = 1 \dots n$:

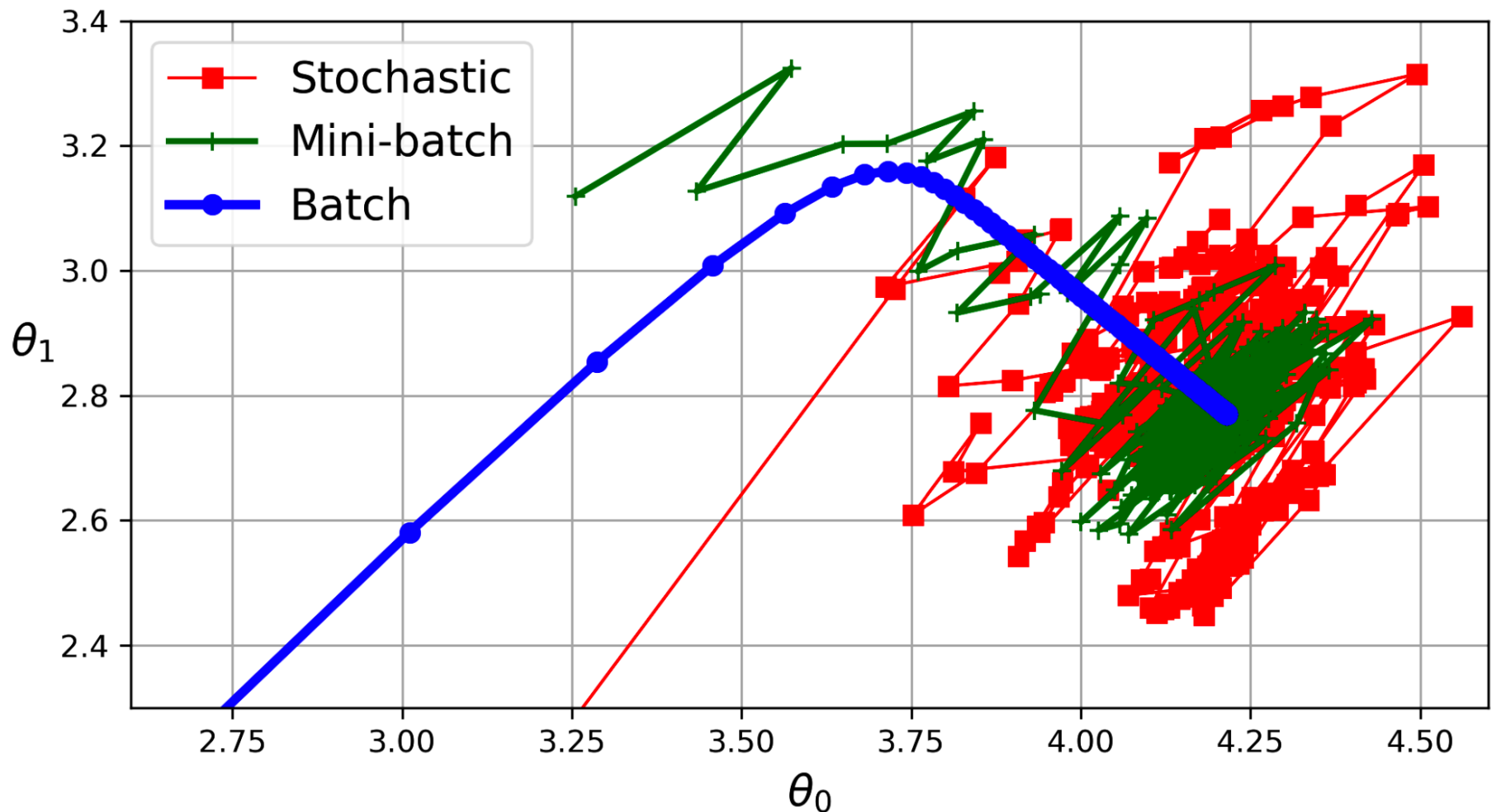
$\mathbf{w} \leftarrow \mathbf{w} - \text{alpha}(\text{derivative of } J(\mathbf{w}) \text{ wrt } x_i)$

 store $J(\mathbf{w})$

Gradient descent variations

- **Batch gradient descent**
 - Go over *all* training data before making a weight update
- **Stochastic gradient descent**
 - Shuffle data and make a weight update after *each training example*
- **Mini-batch gradient descent**
 - Update after a “*mini-batch*” of training examples (i.e. 50)
 - Vocab word: **epoch** (one pass through the data)

Comparison of gradient descent paths



Outline for Feb 20

- Stochastic gradient descent
- Maximizing likelihood functions
- Logistic regression likelihood
- Extending logistic regression to multi-class classification (softmax)
- Regularization

Logistic Regression

$$h_{\vec{w}}(\vec{x}) = p(y=1 | \vec{x})$$

$$y \in \{0, 1\} = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$$

$$\vec{w} \cdot \vec{x} = w_0 + w_1 x_1 + w_2 x_2 \dots w_p x_p$$

↑
"fake 1"

Big Q: what is \vec{w} ?

Likelihood function

n
training
examples

$$L(\vec{w}) = \prod_{i=1}^n \underbrace{h_{\vec{w}}(\vec{x}_i)}_{\text{prob of label 1}}^{y_i} \underbrace{(1 - h_{\vec{w}}(\vec{x}_i))}_{\text{prob of label 0}}^{1 - y_i}$$

product

maximize likelihood

Aside to maximum likelihood estimators (MLE)

coin flip

flips = n

$$\begin{aligned}\text{prob}(1) &= p \\ \text{prob}(0) &= 1-p\end{aligned}$$

$$\vec{y} = [0, 0, 1, 1, 0, 1, 0, 1, 0, 0]$$

$n=10$

$$\begin{aligned}L(p) &= (1-p)(1-p)p p(1-p)p \cdot (1-p) \\ &= p^4(1-p)^6\end{aligned}$$

$$L(p) = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i} = p^{\#1s} (1-p)^{\#0s}$$

$$\#1s = \bar{y}n = \frac{4}{10} \cdot 10$$

mean/avg

$$\#0s = (1-\bar{y})n$$

Handout 9

$$\log(b^a) = a \log b$$

$$f(x) = \log x$$

$$f'(x) = \frac{1}{x}$$

$$l(p) = \bar{y} n \log p + n(1-\bar{y}) \log(1-p)$$

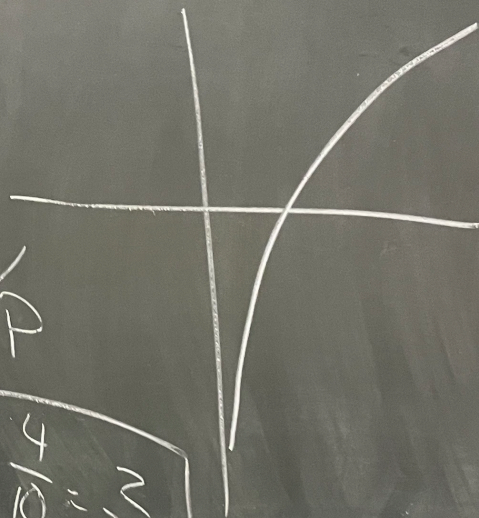
$$l'(p) = \frac{\cancel{n\bar{y}}}{p} - \frac{\cancel{n(1-\bar{y})}}{1-p} = 0$$

$$\frac{\bar{y}}{p} = \frac{1-\bar{y}}{1-p}$$

$$\bar{y} - \bar{y} \cancel{p} = p - \bar{y} \cancel{p}$$

$$\boxed{\hat{p} = \bar{y}}$$

$$\boxed{\hat{p} = \frac{4}{10} = \frac{2}{5}}$$



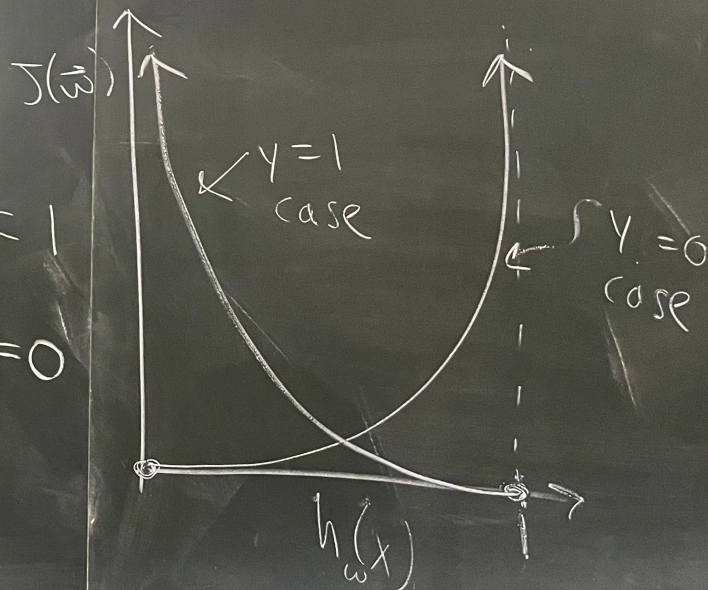
negative log likelihood (minimize)

$$J(\vec{w}) = -\log L(\vec{w})$$

$$= -\sum_{i=1}^n y_i \log h_{\vec{w}}(x_i) + (1-y_i) \log (1-h_{\vec{w}}(x_i))$$

single example

$$J_x(\vec{w}) = \begin{cases} -\log h_{\vec{w}}(\vec{x}), & y=1 \\ -\log(1-h_{\vec{w}}(\vec{x})), & y=0 \end{cases}$$



Stochastic Gradient Descent

Shuffle

for $i=1 \dots n$

$$\vec{w} \leftarrow \vec{w} - \alpha \left[\nabla_{\vec{w}} J_{\vec{x}_i}(\vec{w}) \right] \star$$

$$\nabla J_{\vec{x}_i}(\vec{w}) = \left(\frac{y_i}{h(\vec{x}_i)} - \frac{(1-y_i)}{1-h(\vec{x}_i)} \right) \nabla h_{\vec{w}}(\vec{x}_i)$$

$$g(z) = \frac{1}{1 + e^{-z}}, \quad g'(z) = g(z)(1-g(z))$$

$= (1 + e^{-z})^{-1}$

exercise!

$$\begin{aligned} & \rightarrow \left(\frac{y_i}{h(\vec{x}_i)} - \frac{1-y_i}{1-h(\vec{x}_i)} \right) h(\vec{x}_i)(1-h(\vec{x}_i)) \vec{x}_i \\ & = \left(-y_i + y_i \cancel{h(\vec{x}_i)} + h(\vec{x}_i) - y_i \cancel{h(\vec{x}_i)} \right) \vec{x}_i \\ & = \left(h(\vec{x}_i) - y_i \right) \vec{x}_i \end{aligned}$$

\uparrow pred \uparrow truth

$\left[\nabla J_{\vec{x}_i}(\vec{w}) \right]$

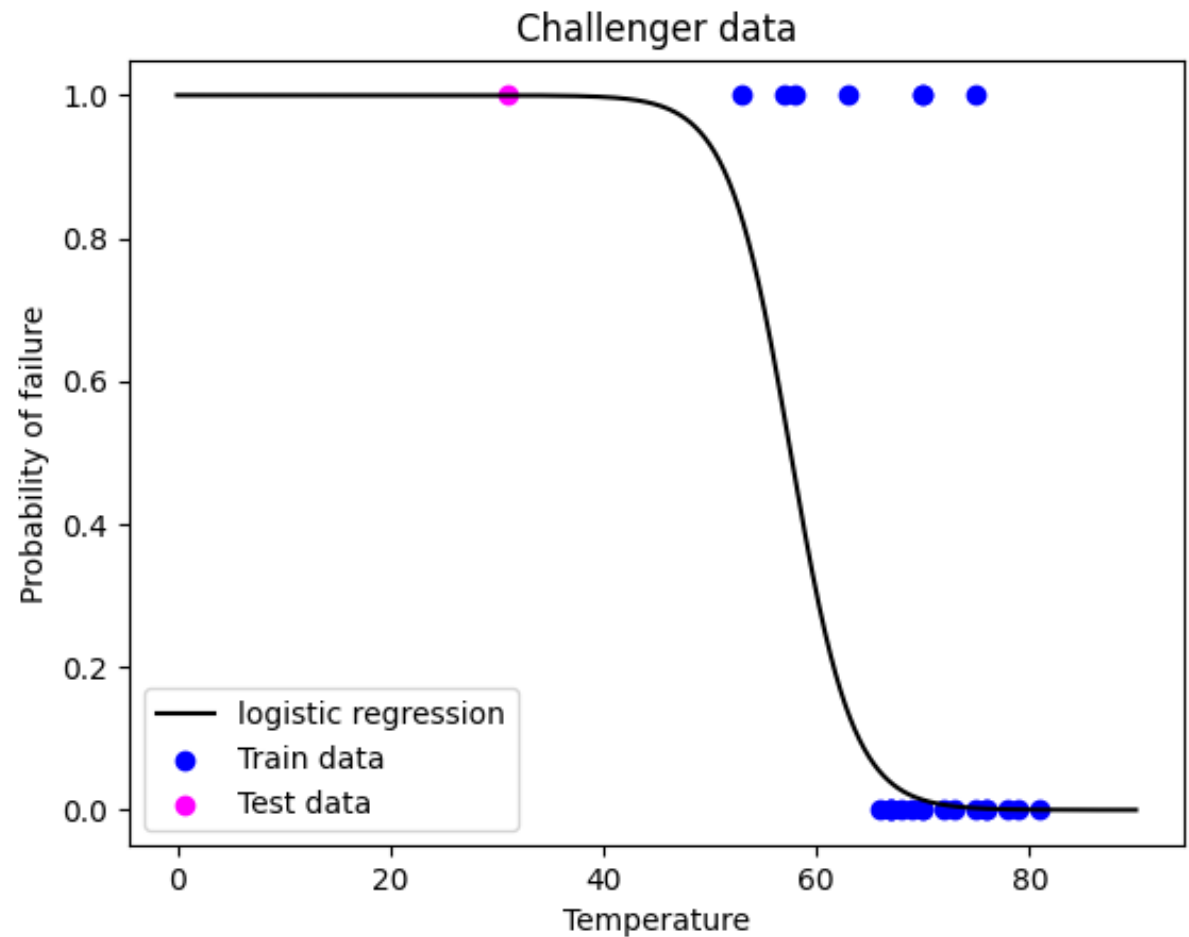
Challenger Explosion Data



Image: NASA

1	Date	Temperature	Damage Incident
2	04/12/1981	66	0
3	11/12/1981	70	1
4	3/22/82	69	0
5	6/27/82	80	NA
6	01/11/1982	68	0
7	04/04/1983	67	0
8	6/18/83	72	0
9	8/30/83	73	0
10	11/28/83	70	0
11	02/03/1984	57	1
⋮			
23	10/30/85	75	1
24	11/26/85	76	0
25	01/12/1986	58	1
26	1/28/86	31	Challenger Accident

Challenger Explosion Data



Logistic Function

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

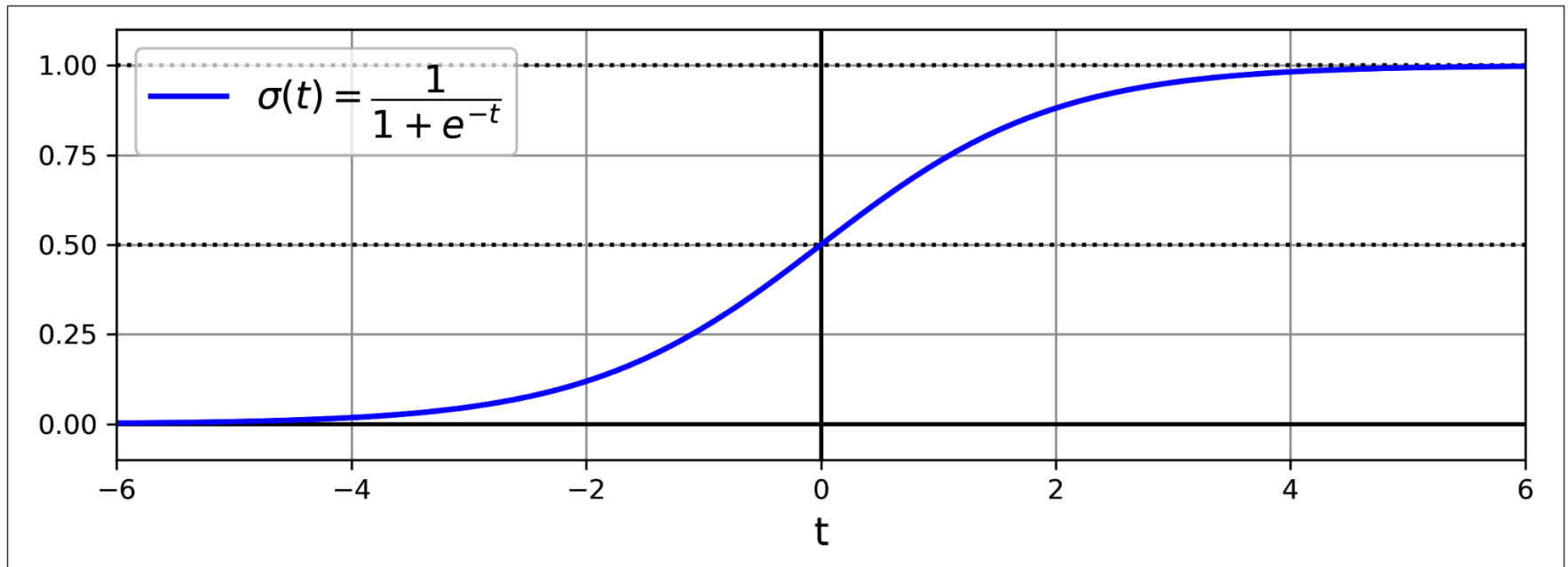


Figure 4-21. Logistic function

Logistic Regression creates a linear decision boundary!

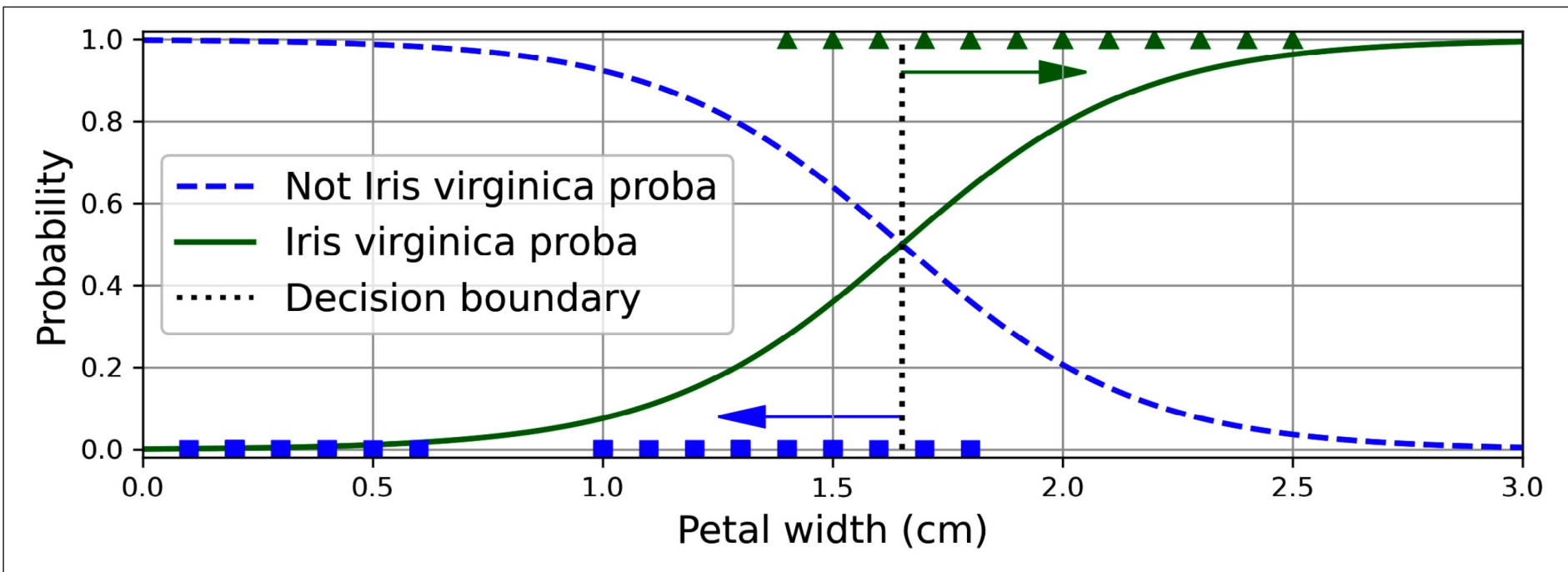


Figure 4-23. Estimated probabilities and decision boundary

Logistic Regression model

- Equivalent to $p(y=1 \mid \mathbf{x})$ from Naïve Bayes

$$\hat{p} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

- Idea of threshold is still the same!

$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5 \\ 1 & \text{if } \hat{p} \geq 0.5 \end{cases}$$

Logistic Regression Cost

- Cost function for single example

$$c(\boldsymbol{\theta}) = \begin{cases} -\log(\hat{p}) & \text{if } y = 1 \\ -\log(1 - \hat{p}) & \text{if } y = 0 \end{cases}$$

- Overall logistic regression cost function

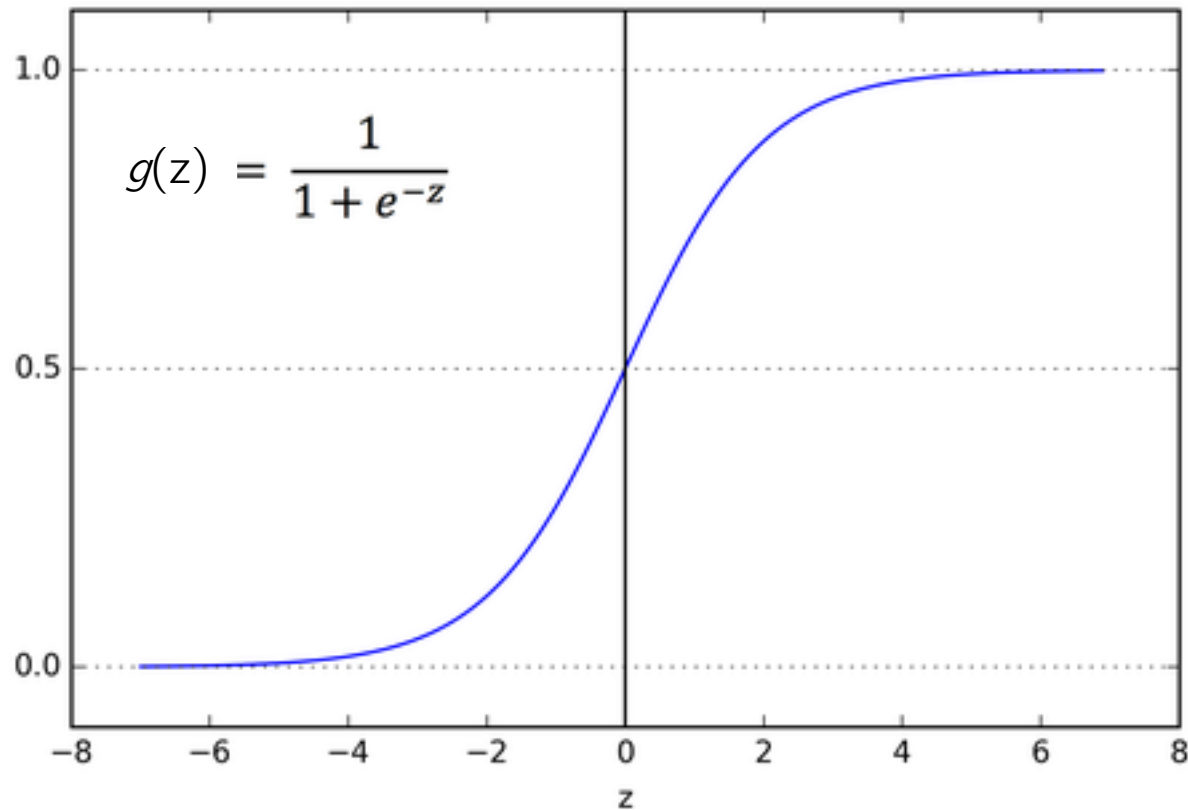
$$J(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right]$$

3 important pieces to SGD

- Hypothesis function (prediction)

$$h_w(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-w \cdot \mathbf{x}}}$$

Logistic (sigmoid) function



3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

3 important pieces to SGD

- Hypothesis function (prediction)

$$h_{\mathbf{w}}(\mathbf{x}) = p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- Cost function (want to minimize)

$$J(\mathbf{w}) = - \sum_{i=1}^n y_i \log h_{\mathbf{w}}(\mathbf{x}_i) + (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient of cost wrt single data point \mathbf{x}_i

$$\nabla J_{\mathbf{x}_i}(\mathbf{w}) = (h_{\mathbf{w}}(\mathbf{x}_i) - y_i) \mathbf{x}_i$$

Outline for Feb 20

- Stochastic gradient descent
- Maximizing likelihood functions
- Logistic regression likelihood
- Extending logistic regression to multi-class classification (softmax)
- Regularization

Multi-class logistic Regression

- political parties
- blood groups

2 classes : $h(\vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} = \frac{e^{\vec{w} \cdot \vec{x}}}{e^{\vec{w} \cdot \vec{x}} + 1}$

$\underbrace{e^{\vec{w} \cdot \vec{x}}}_{\text{weight on class 1}}$
 $+$
 $\underbrace{1}_{\text{weight on class 0}}$

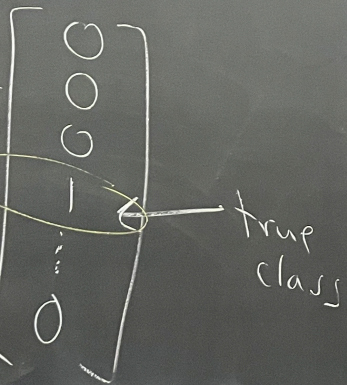
K classes

$$\underbrace{\hat{\vec{y}}}_{\text{prediction}} = \begin{bmatrix} p(y=1|\vec{x}) \\ p(y=2|\vec{x}) \\ \vdots \\ p(y=K|\vec{x}) \end{bmatrix}$$

$\underbrace{\vec{y}}_{\text{true}} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$

→ true class

$$\begin{aligned}
 &P(y=1|\vec{x}) \\
 &P(y=2|\vec{x}) \\
 &\vdots \\
 &P(y=K|\vec{x})
 \end{aligned}$$



$$\frac{1}{\sum_{k=1}^K e^{\vec{w}^{(k)} \cdot \vec{x}}}$$

Diagram illustrating the softmax function. A large bracket groups the terms $e^{\vec{w}^{(1)} \cdot \vec{x}}$, $e^{\vec{w}^{(2)} \cdot \vec{x}}$, and $e^{\vec{w}^{(K)} \cdot \vec{x}}$. Arrows point from these terms to labels s_1 , s_2 , and s_K respectively. Below the bracket, the text "Sum to 1" is written.



$$y_i \log(h(\vec{x}_i)) + (1-y_i) \log(1-h(\vec{x}_i))$$

Diagram illustrating the binary cross-entropy loss function. The equation is written as $y_i \log(h(\vec{x}_i)) + (1-y_i) \log(1-h(\vec{x}_i))$. A box labeled $K=5$ is shown to the right of the equation.