

# CS480/680: Introduction to Machine Learning

## Lec 03: Logistic Regression

Yaoliang Yu



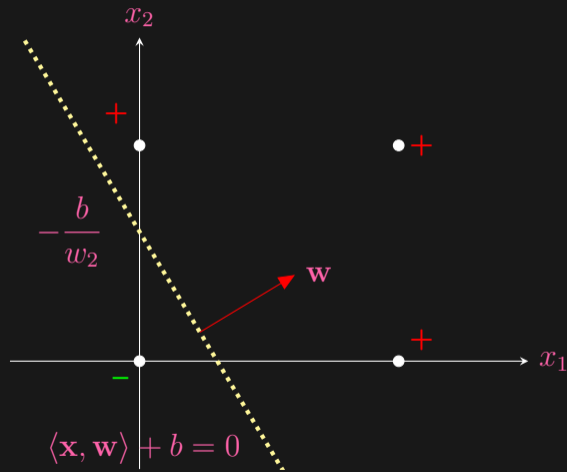
UNIVERSITY OF  
**WATERLOO**

FACULTY OF MATHEMATICS  
DAVID R. CHERITON SCHOOL  
OF COMPUTER SCIENCE

May 15, 2024

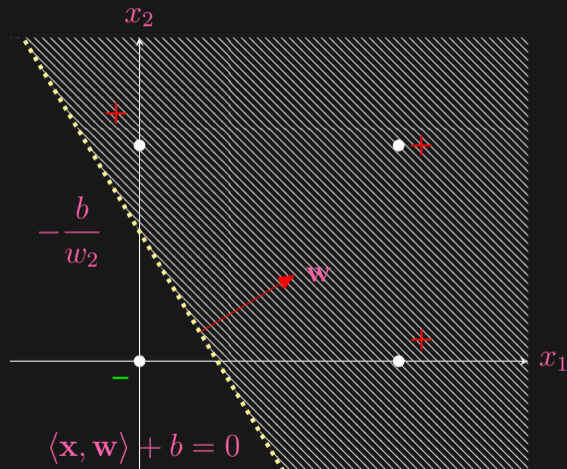
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
- Better idea: learn confidence directly



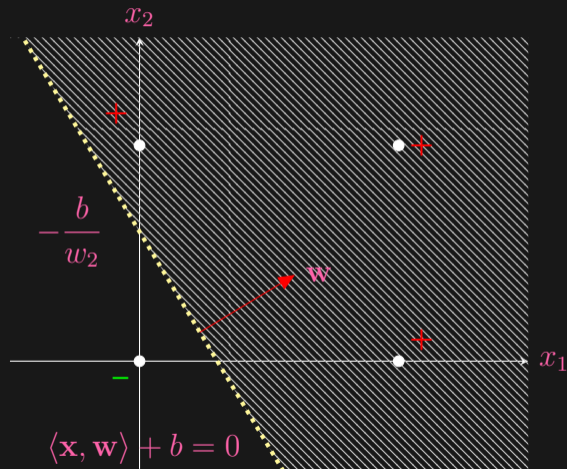
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
- Better idea: learn confidence directly



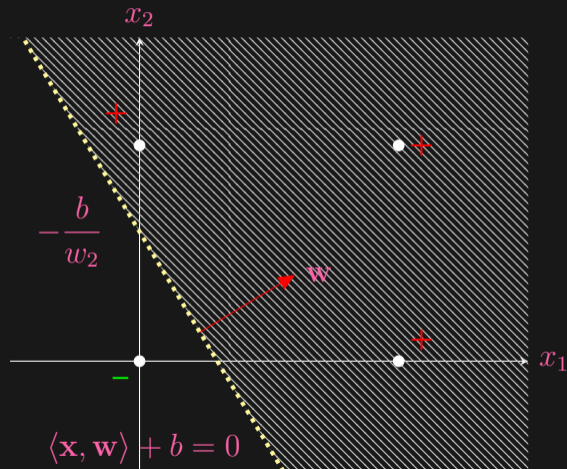
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
- Better idea: learn confidence directly



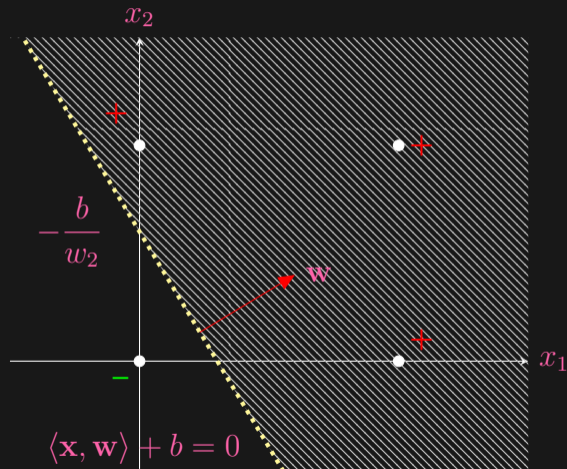
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
  - in fact was used in multi-class perceptron
  - real-valued: hard to interpret
  - many ways to transform into  $[0, 1]$
- Better idea: learn confidence directly



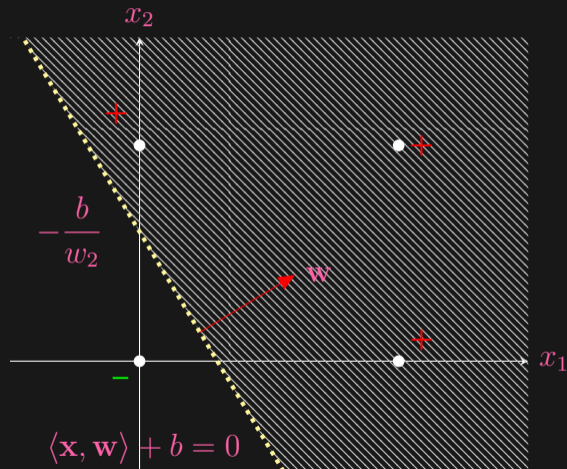
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
  - in fact was used in multi-class perceptron
  - real-valued: hard to interpret
  - many ways to transform into  $[0, 1]$
- Better idea: learn confidence directly



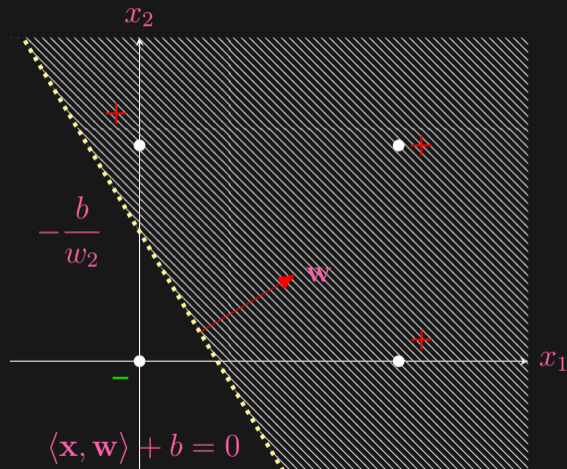
# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
  - in fact was used in multi-class perceptron
  - real-valued: hard to interpret
  - many ways to transform into  $[0, 1]$
- Better idea: learn confidence directly



# Predicting with Confidence

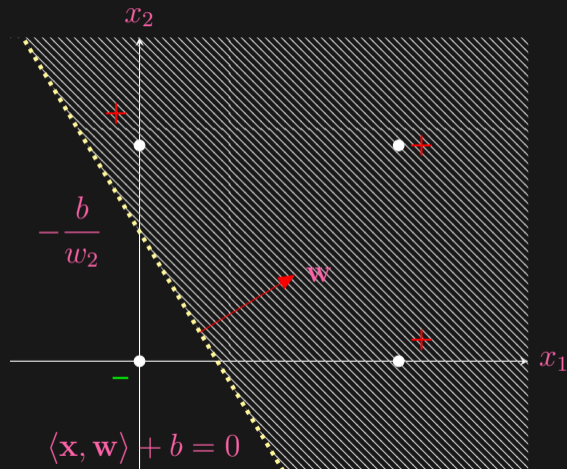
- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
  - in fact was used in multi-class perceptron
  - real-valued: hard to interpret
  - many ways to transform into  $[0, 1]$
- Better idea: learn confidence directly





# Predicting with Confidence

- Recall that  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- How confident we are about the prediction  $\hat{y}$ ?
- Can use  $|\langle \mathbf{x}, \mathbf{w} \rangle|$  as an indication
  - in fact was used in multi-class perceptron
  - real-valued: hard to interpret
  - many ways to transform into  $[0, 1]$
- **Better** idea: learn confidence directly



# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- Scoring function:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- Scoring rule:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) properness (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- Entropy:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- Scoring function:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- Scoring rule:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) properness (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- Entropy:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- Scoring function:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- Scoring rule:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) properness (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- Entropy:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- Scoring function:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- Scoring rule:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) properness (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- Entropy:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- **Scoring function:**  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- **Scoring rule:**  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) **properness** (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- **Entropy:**  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- **Scoring function:**  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- **Scoring rule:**  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) **properness** (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- **Entropy:**  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- **Scoring function**:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- **Scoring rule**:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) **properness** (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- **Entropy**:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$



# Confidence Game

- $Y_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(q)$  for some  $q \in [0, 1]$ 
  - e.g., probability of raining tomorrow
- How to evaluate a probabilistic forecast  $\hat{p}$ ?
- **Scoring function**:  $s : \mathcal{Y} \times [0, 1] \rightarrow \mathbb{R}$ ,  $s(y, p)$  scores the “fitness”
- **Scoring rule**:  $\mathbb{S} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ ,  $\mathbb{S}(q, p) := \mathbb{E}_{Y \sim \text{Bernoulli}(q)}[s(Y, p)]$
- (Strict) **properness** (truthfulness):  $q = \operatorname{argmin}_p \mathbb{S}(q, p)$
- **Entropy**:  $\mathbb{H}(q) := \min_p \mathbb{S}(q, p)$ , under properness,  $\mathbb{H}(q) = \mathbb{S}(q, q)$

# Logarithmic Loss

$$s(\mathbf{y}, p) := -\mathbf{y} \log p - (1 - \mathbf{y}) \log(1 - p)$$

$$\mathbb{S}(q, p) := -q \log p - (1 - q) \log(1 - p)$$

$$\mathbb{H}(q) := -q \log q - (1 - q) \log(1 - q)$$

- Indeed a proper scoring rule (could take  $\infty$  value)
- The resulting entropy is exactly **Shannon's entropy**
- **KL divergence**:  $\text{KL}(q, p) := \mathbb{S}(q, p) - \mathbb{H}(q) \geq 0$ , with equality iff  $q = p$

# Logarithmic Loss

$$s(\mathbf{y}, p) := -\mathbf{y} \log p - (1 - \mathbf{y}) \log(1 - p)$$

$$\mathbb{S}(q, p) := -q \log p - (1 - q) \log(1 - p)$$

$$\mathbb{H}(q) := -q \log q - (1 - q) \log(1 - q)$$

- Indeed a proper scoring rule (could take  $\infty$  value)
- The resulting entropy is exactly Shannon's entropy
- KL divergence:  $\text{KL}(q, p) := \mathbb{S}(q, p) - \mathbb{H}(q) \geq 0$ , with equality iff  $q = p$

# Logarithmic Loss

$$s(\mathbf{y}, p) := -\mathbf{y} \log p - (1 - \mathbf{y}) \log(1 - p)$$

$$\mathbb{S}(q, p) := -q \log p - (1 - q) \log(1 - p)$$

$$\mathbb{H}(q) := -q \log q - (1 - q) \log(1 - q)$$

- Indeed a proper scoring rule (could take  $\infty$  value)
- The resulting entropy is exactly **Shannon's entropy**
- **KL divergence**:  $\text{KL}(q, p) := \mathbb{S}(q, p) - \mathbb{H}(q) \geq 0$ , with equality iff  $q = p$

# Logarithmic Loss

$$s(\mathbf{y}, p) := -\mathbf{y} \log p - (1 - \mathbf{y}) \log(1 - p)$$

$$\mathbb{S}(q, p) := -q \log p - (1 - q) \log(1 - p)$$

$$\mathbb{H}(q) := -q \log q - (1 - q) \log(1 - q)$$

- Indeed a proper scoring rule (could take  $\infty$  value)
- The resulting entropy is exactly **Shannon's entropy**
- **KL divergence**:  $\text{KL}(q, p) := \mathbb{S}(q, p) - \mathbb{H}(q) \geq 0$ , with equality iff  $q = p$

# Introducing $\mathbb{X}$

$$s(\mathbf{y}, p(\mathbf{x})) := -\mathbf{y} \log p(\mathbf{x}) - (1 - \mathbf{y}) \log(1 - p(\mathbf{x}))$$

$$S(q(\mathbf{x}), p(\mathbf{x})) := -q(\mathbf{x}) \log p(\mathbf{x}) - (1 - q(\mathbf{x})) \log(1 - p(\mathbf{x}))$$

$$\mathbb{S}(q, p) := \mathbb{E}[-q(\mathbf{X}) \log p(\mathbf{X}) - (1 - q(\mathbf{X})) \log(1 - p(\mathbf{X}))]$$

- $Y|X = \mathbf{x} \sim \text{Bernoulli}(q(\mathbf{x}))$
- Observe that  $S(q(\mathbf{x}), p(\mathbf{x})) = \mathbb{E}_{Y|X=\mathbf{x}}[s(Y, p(\mathbf{x}))]$ ,  $\mathbb{S}(q, p) := \mathbb{E}[s(Y, p(\mathbf{X}))]$
- Parameterizing the probabilistic forecast, e.g.  $p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- Minimum score estimation:

$$\min_{\mathbf{w}} \hat{\mathbb{E}}[s(Y, p(\mathbf{X}; \mathbf{w}))]$$

# Introducing $\mathbb{X}$

$$s(\mathbf{y}, p(\mathbf{x})) := -\mathbf{y} \log p(\mathbf{x}) - (1 - \mathbf{y}) \log(1 - p(\mathbf{x}))$$

$$S(q(\mathbf{x}), p(\mathbf{x})) := -q(\mathbf{x}) \log p(\mathbf{x}) - (1 - q(\mathbf{x})) \log(1 - p(\mathbf{x}))$$

$$\mathbb{S}(q, p) := \mathbb{E}[-q(\mathbf{X}) \log p(\mathbf{X}) - (1 - q(\mathbf{X})) \log(1 - p(\mathbf{X}))]$$

- $\mathbf{Y}|\mathbf{X} = \mathbf{x} \sim \text{Bernoulli}(q(\mathbf{x}))$
- Observe that  $S(q(\mathbf{x}), p(\mathbf{x})) = \mathbb{E}_{\mathbf{Y}|\mathbf{X}=\mathbf{x}}[s(\mathbf{Y}, p(\mathbf{x}))]$ ,  $\mathbb{S}(q, p) := \mathbb{E}[s(\mathbf{Y}, p(\mathbf{X}))]$
- Parameterizing the probabilistic forecast, e.g.  $p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- Minimum score estimation:

$$\min_{\mathbf{w}} \hat{\mathbb{E}}[s(\mathbf{Y}, p(\mathbf{X}; \mathbf{w}))]$$

# Introducing $\mathbb{X}$

$$s(\mathbf{y}, p(\mathbf{x})) := -\mathbf{y} \log p(\mathbf{x}) - (1 - \mathbf{y}) \log(1 - p(\mathbf{x}))$$

$$S(q(\mathbf{x}), p(\mathbf{x})) := -q(\mathbf{x}) \log p(\mathbf{x}) - (1 - q(\mathbf{x})) \log(1 - p(\mathbf{x}))$$

$$\mathbb{S}(q, p) := \mathbb{E}[-q(\mathbf{X}) \log p(\mathbf{X}) - (1 - q(\mathbf{X})) \log(1 - p(\mathbf{X}))]$$

- $\mathbf{Y} | \mathbf{X} = \mathbf{x} \sim \text{Bernoulli}(q(\mathbf{x}))$
- Observe that  $S(q(\mathbf{x}), p(\mathbf{x})) = \mathbb{E}_{\mathbf{Y} | \mathbf{X} = \mathbf{x}} s(\mathbf{Y}, p(\mathbf{x}))$ ,  $\mathbb{S}(q, p) := \mathbb{E}[s(\mathbf{Y}, p(\mathbf{X}))]$
- Parameterizing the probabilistic forecast, e.g.  $p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- Minimum score estimation:

$$\min_{\mathbf{w}} \hat{\mathbb{E}}[s(\mathbf{Y}, p(\mathbf{X}; \mathbf{w}))]$$



# Introducing $\mathbb{X}$

$$s(\mathbf{y}, p(\mathbf{x})) := -\mathbf{y} \log p(\mathbf{x}) - (1 - \mathbf{y}) \log(1 - p(\mathbf{x}))$$

$$S(q(\mathbf{x}), p(\mathbf{x})) := -q(\mathbf{x}) \log p(\mathbf{x}) - (1 - q(\mathbf{x})) \log(1 - p(\mathbf{x}))$$

$$\mathbb{S}(q, p) := \mathbb{E}[-q(\mathbf{X}) \log p(\mathbf{X}) - (1 - q(\mathbf{X})) \log(1 - p(\mathbf{X}))]$$

- $Y|X = \mathbf{x} \sim \text{Bernoulli}(q(\mathbf{x}))$
- Observe that  $S(q(\mathbf{x}), p(\mathbf{x})) = \mathbb{E}_{Y|X=\mathbf{x}} s(Y, p(\mathbf{x}))$ ,  $\mathbb{S}(q, p) := \mathbb{E}[s(Y, p(\mathbf{X}))]$
- Parameterizing the probabilistic forecast, e.g.  $p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- Minimum score estimation:

$$\min_{\mathbf{w}} \hat{\mathbb{E}}[s(Y, p(\mathbf{X}; \mathbf{w}))]$$

# Introducing $\mathbb{X}$

$$s(\mathbf{y}, p(\mathbf{x})) := -\mathbf{y} \log p(\mathbf{x}) - (1 - \mathbf{y}) \log(1 - p(\mathbf{x}))$$

$$S(q(\mathbf{x}), p(\mathbf{x})) := -q(\mathbf{x}) \log p(\mathbf{x}) - (1 - q(\mathbf{x})) \log(1 - p(\mathbf{x}))$$

$$\mathbb{S}(q, p) := \mathbb{E}[-q(\mathbf{X}) \log p(\mathbf{X}) - (1 - q(\mathbf{X})) \log(1 - p(\mathbf{X}))]$$

- $\mathbf{Y} | \mathbf{X} = \mathbf{x} \sim \text{Bernoulli}(q(\mathbf{x}))$
- Observe that  $S(q(\mathbf{x}), p(\mathbf{x})) = \mathbb{E}_{\mathbf{Y} | \mathbf{X} = \mathbf{x}} s(\mathbf{Y}, p(\mathbf{x}))$ ,  $\mathbb{S}(q, p) := \mathbb{E}[s(\mathbf{Y}, p(\mathbf{X}))]$
- Parameterizing the probabilistic forecast, e.g.  $p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle)$
- Minimum score estimation:

$$\min_{\mathbf{w}} \hat{\mathbb{E}}[s(\mathbf{Y}, p(\mathbf{X}; \mathbf{w}))]$$

# Max Conditional Likelihood

- Model postulates  $Y|X = \mathbf{x} \sim \text{Bernoulli}(p(\mathbf{x}; \mathbf{w}))$ , i.e.  $\Pr(Y = 1|X = \mathbf{x}) = p(\mathbf{x}; \mathbf{w})$
- Given  $(\mathbf{X}_i, y_i), i = 1, \dots, n$ , assume independence:

$$\begin{aligned}\Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n) &= \prod_{i=1}^n \Pr(Y_i = y_i | X_i = \mathbf{x}_i) \\ &= \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}\end{aligned}$$

- Maximizing the conditional log-likelihood:

$$\max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

# Max Conditional Likelihood

- Model postulates  $Y|X = \mathbf{x} \sim \text{Bernoulli}(p(\mathbf{x}; \mathbf{w}))$ , i.e.  $\Pr(Y = 1|X = \mathbf{x}) = p(\mathbf{x}; \mathbf{w})$
- Given  $(X_i, y_i), i = 1, \dots, n$ , assume independence:

$$\begin{aligned}\Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n) &= \prod_{i=1}^n \Pr(Y_i = y_i | X_i = \mathbf{x}_i) \\ &= \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}\end{aligned}$$

- Maximizing the conditional log-likelihood:

$$\max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

# Max Conditional Likelihood

- Model postulates  $Y|X = \mathbf{x} \sim \text{Bernoulli}(p(\mathbf{x}; \mathbf{w}))$ , i.e.  $\Pr(Y = 1|X = \mathbf{x}) = p(\mathbf{x}; \mathbf{w})$
- Given  $(\mathbf{X}_i, y_i), i = 1, \dots, n$ , assume independence:

$$\begin{aligned}\Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n) &= \prod_{i=1}^n \Pr(Y_i = y_i | X_i = \mathbf{x}_i) \\ &= \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}\end{aligned}$$

- Maximizing the conditional log-likelihood:

$$\max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

# Max Conditional Likelihood

- Model postulates  $Y|X = \mathbf{x} \sim \text{Bernoulli}(p(\mathbf{x}; \mathbf{w}))$ , i.e.  $\Pr(Y = 1|X = \mathbf{x}) = p(\mathbf{x}; \mathbf{w})$
- Given  $(\mathbf{X}_i, y_i), i = 1, \dots, n$ , assume independence:

$$\begin{aligned}\Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = \mathbf{x}_1, \dots, X_n = \mathbf{x}_n) &= \prod_{i=1}^n \Pr(Y_i = y_i | X_i = \mathbf{x}_i) \\ &= \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}\end{aligned}$$

- Maximizing the conditional log-likelihood:

$$\max_{\mathbf{w}} \prod_{i=1}^n [p(\mathbf{x}_i; \mathbf{w})]^{y_i} [1 - p(\mathbf{x}_i; \mathbf{w})]^{1-y_i}$$

# Two Extremes

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{w})$ ?
- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{x})$ ?

# Two Extremes

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{w})$ ?
  - i.e. use the same confidence  $p$  for every data point
- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{x})$ ?



# Two Extremes

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{w})$ ?
  - i.e. use the same confidence  $p$  for every data point
- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{x})$ ?

# Two Extremes

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{w})$ ?
  - i.e. use the same confidence  $p$  for every data point
- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{x})$ ?
  - i.e. every data point uses its own confidence  $p(\mathbf{x})$

# Two Extremes

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{w})$ ?
  - i.e. use the same confidence  $p$  for every data point
- What is the solution if  $p(\mathbf{x}; \mathbf{w}) = p(\mathbf{x})$ ?
  - i.e. every data point uses its own confidence  $p(\mathbf{x})$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?

- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$

- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$ 
  - i.e., the odds ratio is an affine function
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$



# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$ 
  - i.e., the odds ratio is an affine function
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# The Logit Transform

- $p(\mathbf{x}; \mathbf{w}) : \mathcal{X} \rightarrow [0, 1]$ , how to parameterize using  $\mathbf{w}$ ?
  - $p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
  - $\log p(\mathbf{x}; \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$ ?
- Logit transform:  $\log \frac{p(\mathbf{x}; \mathbf{w})}{1-p(\mathbf{x}; \mathbf{w})} = \langle \mathbf{x}, \mathbf{w} \rangle$ 
  - i.e., the odds ratio is an affine function
- Equivalently, the sigmoid transformation:  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1+\exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

# Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- Plug in the parameterization  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

$$\min_{\mathbf{w}} \sum_{i=1}^n \boxed{\log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle}$$

- Note the label encoding  $y_i \in \{0, 1\}$ ; if instead,  $y_i \in \{\pm 1\}$ , then

$$\min_{\mathbf{w}} \sum_{i=1}^n \underbrace{\log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]}_{\text{logistic loss}}$$

# Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

- Plug in the parameterization  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

$$\min_{\mathbf{w}} \sum_{i=1}^n \boxed{\log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle}$$

- Note the label encoding  $y_i \in \{0, 1\}$ ; if instead,  $y_i \in \{\pm 1\}$ , then

$$\min_{\mathbf{w}} \sum_{i=1}^n \underbrace{\log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]}_{\text{logistic loss}}$$

# Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n -y_i \log[p(\mathbf{x}_i; \mathbf{w})] - (1 - y_i) \log[1 - p(\mathbf{x}_i; \mathbf{w})]$$

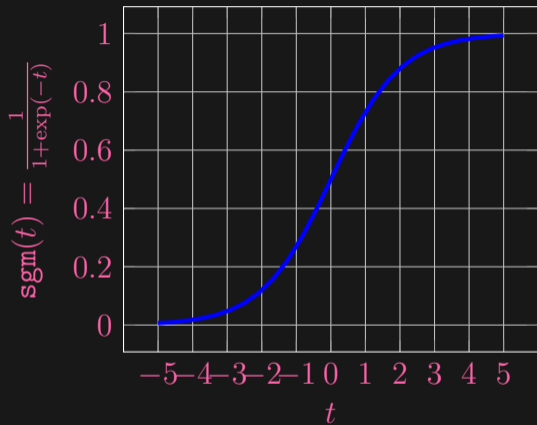
- Plug in the parameterization  $p(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$

$$\min_{\mathbf{w}} \sum_{i=1}^n \boxed{\log[1 + \exp(-\langle \mathbf{x}_i, \mathbf{w} \rangle)] + (1 - y_i) \langle \mathbf{x}_i, \mathbf{w} \rangle}$$

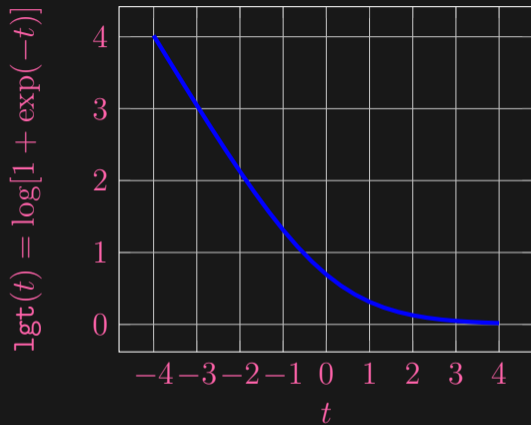
- Note the label encoding  $y_i \in \{0, 1\}$ ; if instead,  $y_i \in \{\pm 1\}$ , then

$$\min_{\mathbf{w}} \sum_{i=1}^n \underbrace{\boxed{\log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]}}_{\text{logistic loss}}$$

### sigmoid function



### logistic loss



D. R. Cox. "The Regression Analysis of Binary Sequences". *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 20, no. 2 (1958), pp. 215–242.

# Prediction

$$p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

- $\hat{y} = 1$  iff  $p(\mathbf{x}; \mathbf{w}) = \Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2}$  iff  $\langle \mathbf{x}, \mathbf{w} \rangle > 0$
- Decision boundary remains to be  $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = 0\}$
- Can predict  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$  as before, but now with confidence  $p(\mathbf{x}; \mathbf{w})$

# Prediction

$$p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

- $\hat{y} = 1$  iff  $p(\mathbf{x}; \mathbf{w}) = \Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2}$  iff  $\langle \mathbf{x}, \mathbf{w} \rangle > 0$
- Decision boundary remains to be  $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = 0\}$
- Can predict  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$  as before, but now with confidence  $p(\mathbf{x}; \mathbf{w})$



# Prediction

$$p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

- $\hat{y} = 1$  iff  $p(\mathbf{x}; \mathbf{w}) = \Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2}$  iff  $\langle \mathbf{x}, \mathbf{w} \rangle > 0$
- Decision boundary remains to be  $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = 0\}$
- Can predict  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$  as before, but now with confidence  $p(\mathbf{x}; \mathbf{w})$

# Prediction

$$p(\mathbf{x}; \mathbf{w}) = \text{sgm}(\langle \mathbf{x}, \mathbf{w} \rangle) = \frac{1}{1 + \exp(-\langle \mathbf{x}, \mathbf{w} \rangle)}$$

- $\hat{y} = 1$  iff  $p(\mathbf{x}; \mathbf{w}) = \Pr(Y = 1|X = \mathbf{x}) > \frac{1}{2}$  iff  $\langle \mathbf{x}, \mathbf{w} \rangle > 0$
- Decision boundary remains to be  $H := \{\mathbf{x} : \langle \mathbf{x}, \mathbf{w} \rangle = 0\}$
- Can predict  $\hat{y} = \text{sign}(\langle \mathbf{x}, \mathbf{w} \rangle)$  as before, but now with confidence  $p(\mathbf{x}; \mathbf{w})$

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1|X = \mathbf{x})$  under the **linear odds ratio assumption**
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!

be lazy

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1 | X = \mathbf{x})$  under the **linear odds ratio assumption**
  - confidence is meaningless if the assumption is way off
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!

be lazy

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1 | X = \mathbf{x})$  under the **linear odds ratio assumption**
  - confidence is meaningless if the assumption is way off
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!

be lazy

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1 | X = \mathbf{x})$  under the **linear odds ratio assumption**
  - confidence is meaningless if the assumption is way off
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!

be lazy

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1 | X = \mathbf{x})$  under the **linear odds ratio assumption**
  - confidence is meaningless if the assumption is way off
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!
  - sufficient but not necessary, be lazy... SVM later

# More than a Classification Algorithm

- Logistic regression estimates the posterior probability  $\eta(\mathbf{x}) := \Pr(Y = 1 | X = \mathbf{x})$  under the **linear odds ratio assumption**
  - confidence is meaningless if the assumption is way off
- Classification itself only requires comparing  $\eta(\mathbf{x})$  with  $\frac{1}{2}$
- Possible to do the comparison without estimating  $\eta(\mathbf{x})$  explicitly!
  - sufficient but not necessary, **be lazy...** SVM later



# Beyond Logistic

$$p(\mathbf{x}; \mathbf{w}) = F(\langle \mathbf{x}, \mathbf{w} \rangle)$$

- $F : \mathbb{R} \rightarrow [0, 1]$ , increasing: any cumulative distribution function (cdf) would do
- Logistic distribution:  $F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x - \mu}{s})}$

# Beyond Logistic

$$p(\mathbf{x}; \mathbf{w}) = F(\langle \mathbf{x}, \mathbf{w} \rangle)$$

- $F : \mathbb{R} \rightarrow [0, 1]$ , increasing: any cumulative distribution function (cdf) would do
- Logistic distribution:  $F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x - \mu}{s})}$

# Beyond Logistic

$$p(\mathbf{x}; \mathbf{w}) = F(\langle \mathbf{x}, \mathbf{w} \rangle)$$

- $F : \mathbb{R} \rightarrow [0, 1]$ , increasing: any cumulative distribution function (cdf) would do
- Logistic distribution:  $F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x - \mu}{s})}$ 
  - with mean  $\mu$  and variance  $s^2 \pi^2 / 3$
  - sigmoid is exactly when  $\mu = 0$  and  $s = 1$

# Beyond Logistic

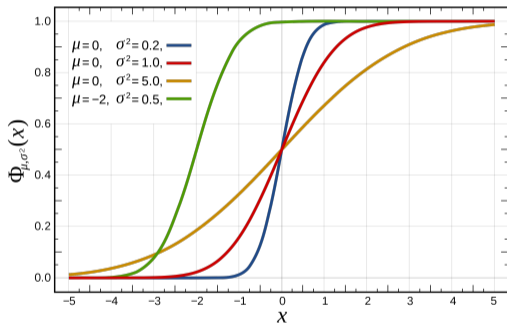
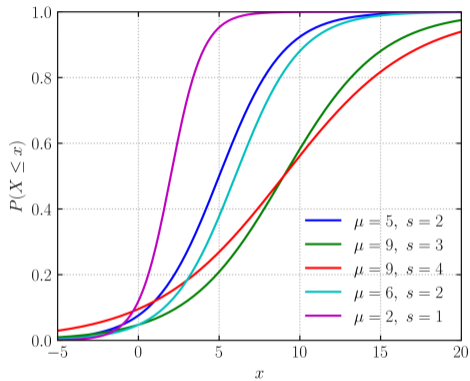
$$p(\mathbf{x}; \mathbf{w}) = F(\langle \mathbf{x}, \mathbf{w} \rangle)$$

- $F : \mathbb{R} \rightarrow [0, 1]$ , increasing: any cumulative distribution function (cdf) would do
- Logistic distribution:  $F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x - \mu}{s})}$ 
  - with mean  $\mu$  and variance  $s^2 \pi^2 / 3$
  - sigmoid is exactly when  $\mu = 0$  and  $s = 1$

# Beyond Logistic

$$p(\mathbf{x}; \mathbf{w}) = F(\langle \mathbf{x}, \mathbf{w} \rangle)$$

- $F : \mathbb{R} \rightarrow [0, 1]$ , increasing: any cumulative distribution function (cdf) would do
- Logistic distribution:  $F(x; \mu, s) = \frac{1}{1 + \exp(-\frac{x - \mu}{s})}$ 
  - with mean  $\mu$  and variance  $s^2 \pi^2 / 3$
  - sigmoid is exactly when  $\mu = 0$  and  $s = 1$



# Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

- Newton's algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot [\nabla^2 f(\mathbf{w})]^{-1} \cdot \nabla f(\mathbf{w})$$

- The gradient  $\nabla f(\mathbf{w}) = \mathbf{X}(\hat{\mathbf{p}} - \frac{\mathbf{y}+1}{2})$ : changing target
- The Hessian  $\nabla^2 f(\mathbf{w}) = \sum_i \hat{p}_i(1 - \hat{p}_i)\mathbf{x}_i\mathbf{x}_i^\top$ : weighted by confidence
- The confidence  $\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)$

# Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

- Newton's algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot [\nabla^2 f(\mathbf{w})]^{-1} \cdot \nabla f(\mathbf{w})$$

- The gradient  $\nabla f(\mathbf{w}) = \mathbf{X}(\hat{\mathbf{p}} - \frac{\mathbf{y}+1}{2})$ : changing target
- The Hessian  $\nabla^2 f(\mathbf{w}) = \sum_i \hat{p}_i(1 - \hat{p}_i)\mathbf{x}_i\mathbf{x}_i^\top$ : weighted by confidence
- The confidence  $\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)$



# Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

- Newton's algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot [\nabla^2 f(\mathbf{w})]^{-1} \cdot \nabla f(\mathbf{w})$$

- The gradient  $\nabla f(\mathbf{w}) = \mathbf{X}(\hat{\mathbf{p}} - \frac{\mathbf{y}+1}{2})$ : changing target
- The Hessian  $\nabla^2 f(\mathbf{w}) = \sum_i \hat{p}_i(1 - \hat{p}_i)\mathbf{x}_i\mathbf{x}_i^\top$ : weighted by confidence
- The confidence  $\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)$

# Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

- Newton's algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot [\nabla^2 f(\mathbf{w})]^{-1} \cdot \nabla f(\mathbf{w})$$

- The gradient  $\nabla f(\mathbf{w}) = \mathbf{X}(\hat{\mathbf{p}} - \frac{\mathbf{y}+1}{2})$ : changing target
- The Hessian  $\nabla^2 f(\mathbf{w}) = \sum_i \hat{p}_i(1 - \hat{p}_i)\mathbf{x}_i\mathbf{x}_i^\top$ : weighted by confidence
- The confidence  $\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)$

# Solving Logistic Regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log[1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)]$$

- Newton's algorithm:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot [\nabla^2 f(\mathbf{w})]^{-1} \cdot \nabla f(\mathbf{w})$$

- The gradient  $\nabla f(\mathbf{w}) = \mathbf{X}(\hat{\mathbf{p}} - \frac{\mathbf{y}+1}{2})$ : changing target
- The Hessian  $\nabla^2 f(\mathbf{w}) = \sum_i \hat{p}_i(1 - \hat{p}_i)\mathbf{x}_i\mathbf{x}_i^\top$ : weighted by confidence
- The confidence  $\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)$

# Linear Regression vs. Logistic Regression

- |   |   |
|---|---|
| <ul style="list-style-type: none"><li>• least-squares: <math>\sum_{i=1}^n (y_i - \hat{y}_i)^2</math></li><li>• prediction: <math>\hat{y}_i = \langle \mathbf{x}_i, \mathbf{w} \rangle</math></li><li>• objective: <math>\ \mathbf{y} - \hat{\mathbf{y}}\ _2^2</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li></ul> | <ul style="list-style-type: none"><li>• cross-entropy: <math>\sum_{i=1}^n -\frac{1+y_i}{2} \log \hat{p}_i - \frac{1-y_i}{2} \log(1-\hat{p}_i)</math></li><li>• prediction: <math>\hat{y}_i = \text{sign}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math>, <math>\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math></li><li>• objective: <math>\text{KL}(\frac{1+\mathbf{y}}{2} \parallel \hat{\mathbf{p}})</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\hat{S}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li></ul> |
|---|---|

- Diagonal weight matrix  $\hat{S} = \text{diag}(\hat{\mathbf{p}} \odot (1 - \hat{\mathbf{p}}))$
- Logistic regression = iteratively weighted linear regression

# Linear Regression vs. Logistic Regression

- |   |  |
|---|--|
| <ul style="list-style-type: none"><li>• least-squares: <math>\sum_{i=1}^n (y_i - \hat{y}_i)^2</math></li><li>• prediction: <math>\hat{y}_i = \langle \mathbf{x}_i, \mathbf{w} \rangle</math></li><li>• objective: <math>\ \mathbf{y} - \hat{\mathbf{y}}\ _2^2</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li></ul> | <ul style="list-style-type: none"><li>• cross-entropy: <math>\sum_{i=1}^n -\frac{1+y_i}{2} \log \hat{p}_i - \frac{1-y_i}{2} \log(1-\hat{p}_i)</math></li><li>• prediction: <math>\hat{y}_i = \text{sign}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math>, <math>\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math></li><li>• objective: <math>\text{KL}(\frac{1+\mathbf{y}}{2} \parallel \hat{\mathbf{p}})</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\hat{\mathbf{S}}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li></ul> |
|---|--|

- Diagonal weight matrix  $\hat{\mathbf{S}} = \text{diag}(\hat{\mathbf{p}} \odot (1 - \hat{\mathbf{p}}))$
- Logistic regression = iteratively weighted linear regression

# Linear Regression vs. Logistic Regression

- |   |   |
|---|---|
| <ul style="list-style-type: none"><li>• least-squares: <math>\sum_{i=1}^n (y_i - \hat{y}_i)^2</math></li><li>• prediction: <math>\hat{y}_i = \langle \mathbf{x}_i, \mathbf{w} \rangle</math></li><li>• objective: <math>\ \mathbf{y} - \hat{\mathbf{y}}\ _2^2</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{y}} - \mathbf{y})</math></li></ul> | <ul style="list-style-type: none"><li>• cross-entropy: <math>\sum_{i=1}^n -\frac{1+y_i}{2} \log \hat{p}_i - \frac{1-y_i}{2} \log(1-\hat{p}_i)</math></li><li>• prediction: <math>\hat{y}_i = \text{sign}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math>, <math>\hat{p}_i = \text{sgm}(\langle \mathbf{x}_i, \mathbf{w} \rangle)</math></li><li>• objective: <math>\text{KL}(\frac{1+\mathbf{y}}{2} \parallel \hat{\mathbf{p}})</math></li><li>• grad: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li><li>• Newton: <math>\mathbf{w} \leftarrow \mathbf{w} - \eta (\mathbf{X}\hat{S}\mathbf{X}^\top)^{-1} \mathbf{X}(\hat{\mathbf{p}} - \frac{1+\mathbf{y}}{2})</math></li></ul> |
|---|---|

- Diagonal weight matrix  $\hat{S} = \text{diag}(\hat{\mathbf{p}} \odot (1 - \hat{\mathbf{p}}))$
- Logistic regression = **iteratively** weighted linear regression

# More than 2 Classes

- Softmax parameterization:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

- Encode  $y \in \{1, \dots, c\}$
- Minimizing again the logarithmic loss:

$$\min_{\mathbf{w}} \hat{\mathbb{E}} \left[ -\log \frac{\exp(\langle X, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle X, \mathbf{w}_l \rangle)} \right]$$

# More than 2 Classes

- **Softmax** parameterization:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

– nonnegative and sum to 1

- Encode  $y \in \{1, \dots, c\}$
- Minimizing again the logarithmic loss:

$$\min_{\mathbf{w}} \hat{\mathbb{E}} \left[ -\log \frac{\exp(\langle X, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle X, \mathbf{w}_l \rangle)} \right]$$



# More than 2 Classes

- **Softmax** parameterization:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

– nonnegative and sum to 1

- Encode  $y \in \{1, \dots, c\}$
- Minimizing again the logarithmic loss:

$$\min_{\mathbf{w}} \hat{\mathbb{E}} \left[ -\log \frac{\exp(\langle X, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle X, \mathbf{w}_l \rangle)} \right]$$

# More than 2 Classes

- **Softmax** parameterization:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

– nonnegative and sum to 1

- Encode  $\mathbf{y} \in \{1, \dots, c\}$
- Minimizing again the logarithmic loss:

$$\min_{\mathbf{w}} \hat{\mathbb{E}} \left[ -\log \frac{\exp(\langle X, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle X, \mathbf{w}_l \rangle)} \right]$$

# More than 2 Classes

- **Softmax** parameterization:

$$\Pr(Y = k | X = \mathbf{x}; \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_c]) = \frac{\exp(\langle \mathbf{x}, \mathbf{w}_k \rangle)}{\sum_{l=1}^c \exp(\langle \mathbf{x}, \mathbf{w}_l \rangle)}$$

– nonnegative and sum to 1

- Encode  $y \in \{1, \dots, c\}$
- Minimizing again the logarithmic loss:

$$\min_{\mathbf{w}} \hat{\mathbb{E}} \left[ -\log \frac{\exp(\langle X, \mathbf{w}_Y \rangle)}{\sum_{l=1}^c \exp(\langle X, \mathbf{w}_l \rangle)} \right]$$

