

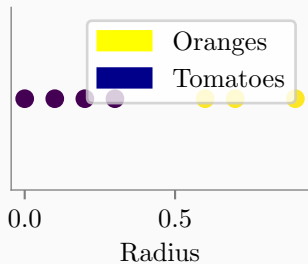
Logistic Regression

Nipun Batra

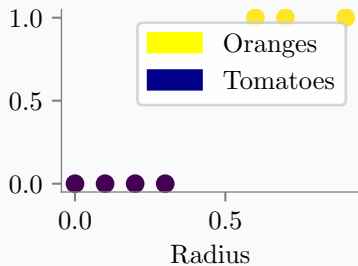
January 2, 2024

IIT Gandhinagar

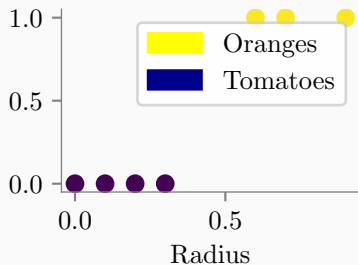
Classification Technique



Classification Technique

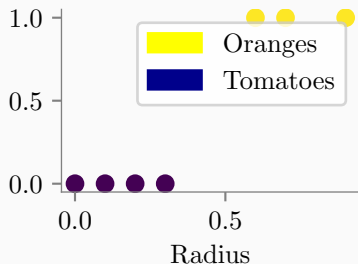


Classification Technique



Aim: $\text{Probability}(\text{Tomatoes} \mid \text{Radius})$? or

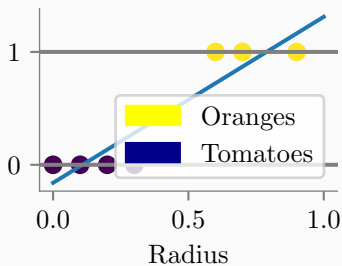
Classification Technique



Aim: Probability(Tomatoes | Radius) ? or

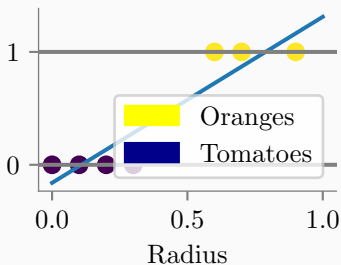
More generally, $P(y = 1|X = x)$?

Idea: Use Linear Regression



$$P(X = \text{Orange} | \text{Radius}) = \theta_0 + \theta_1 \times \text{Radius}$$

Idea: Use Linear Regression



$$P(X = \text{Orange} | \text{Radius}) = \theta_0 + \theta_1 \times \text{Radius}$$

Generally,

$$P(y = 1 | x) = X\theta$$

Idea: Use Linear Regression

Prediction:

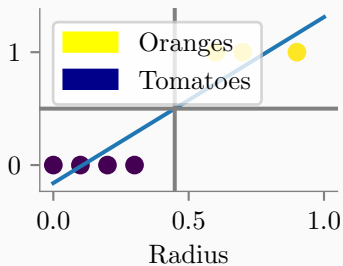
If $\theta_0 + \theta_1 \times \text{Radius} > 0.5 \rightarrow \text{Orange}$
Else $\rightarrow \text{Tomato}$

Problem:

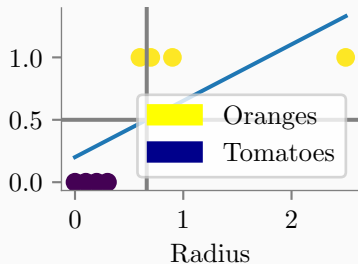
Range of $X\theta$ is $(-\infty, \infty)$

But $P(y = 1 | \dots) \in [0, 1]$

Idea: Use Linear Regression

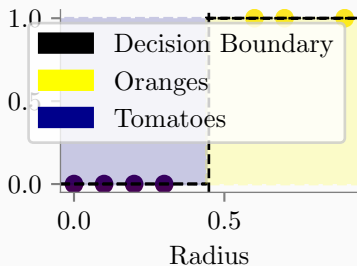


Idea: Use Linear Regression



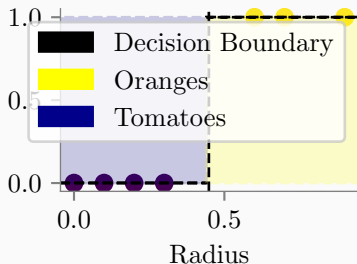
Linear regression for classification gives a poor prediction!

Ideal boundary



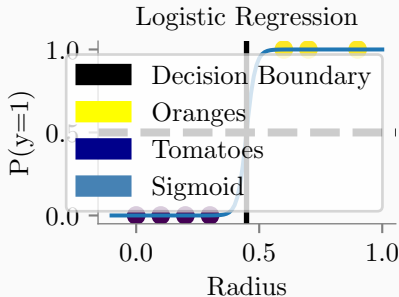
- Have a decision function similar to the above (but not so sharp and discontinuous)

Ideal boundary



- Have a decision function similar to the above (but not so sharp and discontinuous)
- Aim: use linear regression still!

Idea: Use Linear Regression



Question. Can we still use Linear Regression?

Answer. Yes! Transform $\hat{y} \rightarrow [0, 1]$

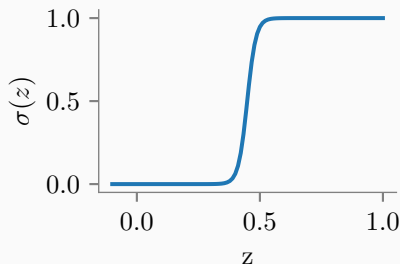
Logistic / Sigmoid Function

$$\hat{y} \in (-\infty, \infty)$$

ϕ = Sigmoid / Logistic Function (σ)

$$\phi(\hat{y}) \in [0, 1]$$

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



Logistic / Sigmoid Function

$$z \rightarrow \infty$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty$$

$$\sigma(z) \rightarrow 0$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty$$

$$\sigma(z) \rightarrow 0$$

$$z = 0$$

Logistic / Sigmoid Function

$$z \rightarrow \infty$$

$$\sigma(z) \rightarrow 1$$

$$z \rightarrow -\infty$$

$$\sigma(z) \rightarrow 0$$

$$z = 0$$

$$\sigma(z) = 0.5$$

Question. Could you use some other transformation (ϕ) of \hat{y} s.t.

$$\phi(\hat{y}) \in [0, 1]$$

Yes! But Logistic Regression works.

$$P(y = 1|X) = \sigma(X\theta) = \frac{1}{1 + e^{-X\theta}}$$

Q. Write $X\theta$ in a more convenient form (as $P(y = 1|X)$, $P(y = 0|X)$)

Logistic / Sigmoid Function

$$P(y = 1|X) = \sigma(X\theta) = \frac{1}{1 + e^{-X\theta}}$$

Q. Write $X\theta$ in a more convenient form (as $P(y = 1|X)$, $P(y = 0|X)$)

Logistic / Sigmoid Function

$$P(y = 1|X) = \sigma(X\theta) = \frac{1}{1 + e^{-X\theta}}$$

Q. Write $X\theta$ in a more convenient form (as $P(y = 1|X)$, $P(y = 0|X)$)

$$P(y = 0|X) = 1 - P(y = 1|X) = 1 - \frac{1}{1 + e^{-X\theta}} = \frac{e^{-X\theta}}{1 + e^{-X\theta}}$$

Logistic / Sigmoid Function

$$P(y = 1|X) = \sigma(X\theta) = \frac{1}{1 + e^{-X\theta}}$$

Q. Write $X\theta$ in a more convenient form (as $P(y = 1|X)$,
 $P(y = 0|X)$)

$$P(y = 0|X) = 1 - P(y = 1|X) = 1 - \frac{1}{1 + e^{-X\theta}} = \frac{e^{-X\theta}}{1 + e^{-X\theta}}$$

$$\therefore \frac{P(y = 1|X)}{1 - P(y = 1|X)} = e^{X\theta} \implies X\theta = \log \frac{P(y = 1|X)}{1 - P(y = 1|X)}$$

Odds (Used in betting)

$$\frac{P(win)}{P(loss)}$$

Here,

$$Odds = \frac{P(y = 1)}{P(y = 0)}$$

$$\log\text{-odds} = \log \frac{P(y=1)}{P(y=0)} = X\theta$$

Q. What is decision boundary for Logistic Regression?

Q. What is decision boundary for Logistic Regression?

Decision Boundary: $P(y = 1|X) = P(y = 0|X)$

$$\text{or } \frac{1}{1+e^{-X\theta}} = \frac{e^{-X\theta}}{1+e^{-X\theta}}$$

$$\text{or } e^{X\theta} = 1$$

$$\text{or } X\theta = 0$$

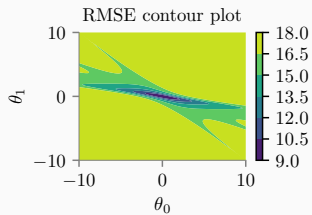
Could we use cost function as:

$$J(\theta) = \sum (y_i - \hat{y}_i)^2$$

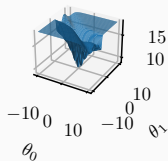
$$\hat{y}_i = \sigma(X\theta)$$

Answer: No (Non-Convex)
(See Jupyter Notebook)

Cost function convexity



RMSE surface plot



$$\text{Likelihood} = P(D|\theta)$$

$$P(y|X, \theta) = \prod_{i=1}^n P(y_i|x_i, \theta)$$

where $y = 0$ or 1

Learning Parameters


$$\text{Likelihood} = P(D|\theta)$$

$$\begin{aligned} P(y|X, \theta) &= \prod_{i=1}^n P(y_i|x_i, \theta) \\ &= \prod_{i=1}^n \left\{ \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{y_i} \left\{ 1 - \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{1-y_i} \end{aligned}$$

[Above: Similar to $P(D|\theta)$ for Linear Regression;
Difference Bernoulli instead of Gaussian]

$$\begin{aligned} -\log P(y|X, \theta) &= \text{Negative Log Likelihood} \\ &= \text{Cost function will be minimising} \\ &= J(\theta) \end{aligned}$$

Likelihood Visualisation



`../figures/logistic-regression/logistic-likelihood.pdf`

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).
- What is $p(H)$?

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).
- What is $p(H)$?
- We might think it to be: $4/10 = 0.4$. But why?

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).
- What is $p(H)$?
- We might think it to be: $4/10 = 0.4$. But why?
- Answer 1: Probability defined as a measure of long running frequencies

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).
- What is $p(H)$?
- We might think it to be: $4/10 = 0.4$. But why?
- Answer 1: Probability defined as a measure of long running frequencies
- Answer 2: What is likelihood of seeing the above sequence when the $p(\text{Head})=\theta$?

Aside on Bernoulli Likelihood

- Assume you have a coin and flip it ten times and get (H, H, T, T, T, H, H, T, T, T).
- What is $p(H)$?
- We might think it to be: $4/10 = 0.4$. But why?
- Answer 1: Probability defined as a measure of long running frequencies
- Answer 2: What is likelihood of seeing the above sequence when the $p(\text{Head})=\theta$?
- Idea find MLE estimate for θ

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$
- What is $P(D_1, D_2, \dots, D_n|\theta)$?

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$
- What is $P(D_1, D_2, \dots, D_n|\theta)$?
- $P(D_1, D_2, \dots, D_n|\theta) = P(D_1|\theta)P(D_2|\theta)\dots P(D_n|\theta)$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$
- What is $P(D_1, D_2, \dots, D_n|\theta)$?
- $P(D_1, D_2, \dots, D_n|\theta) = P(D_1|\theta)P(D_2|\theta)\dots P(D_n|\theta)$
- $P(D_1, D_2, \dots, D_n|\theta) = \theta^{n_h}(1 - \theta)^{n_t}$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$
- What is $P(D_1, D_2, \dots, D_n|\theta)$?
- $P(D_1, D_2, \dots, D_n|\theta) = P(D_1|\theta)P(D_2|\theta)\dots P(D_n|\theta)$
- $P(D_1, D_2, \dots, D_n|\theta) = \theta^{n_h}(1 - \theta)^{n_t}$
- Log-likelihood $= \mathcal{LL}(\theta) = n_h \log(\theta) + n_t \log(1 - \theta)$

Aside on Bernoulli Likelihood

- $p(H) = \theta$ and $p(T) = 1 - \theta$
- What is the PMF for first observation $P(D_1 = x|\theta)$, where $x = 0$ for Tails and $x = 1$ for Heads?
- $P(D_1 = x|\theta) = \theta^x(1 - \theta)^{(1-x)}$
- Verify the above: if $x = 0$ (Tails), $P(D_1 = x|\theta) = 1 - \theta$ and if $x = 1$ (Heads), $P(D_1 = x|\theta) = \theta$
- What is $P(D_1, D_2, \dots, D_n|\theta)$?
- $P(D_1, D_2, \dots, D_n|\theta) = P(D_1|\theta)P(D_2|\theta)\dots P(D_n|\theta)$
- $P(D_1, D_2, \dots, D_n|\theta) = \theta^{n_h}(1 - \theta)^{n_t}$
- Log-likelihood = $\mathcal{LL}(\theta) = n_h \log(\theta) + n_t \log(1 - \theta)$
- $\frac{\partial \mathcal{LL}(\theta)}{\partial \theta} = 0 \implies \frac{n_h}{\theta} + \frac{n_t}{1-\theta} = 0 \implies \theta_{MLE} = \frac{n_h}{n_h + n_t}$

$$J(\theta) = -\log \left\{ \prod_{i=1}^n \left\{ \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{y_i} \left\{ 1 - \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{1-y_i} \right\}$$

$$J(\theta) = -\left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

$$J(\theta) = -\log \left\{ \prod_{i=1}^n \left\{ \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{y_i} \left\{ 1 - \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{1-y_i} \right\}$$
$$J(\theta) = -\left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

This cost function is called cross-entropy.

$$J(\theta) = -\log \left\{ \prod_{i=1}^n \left\{ \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{y_i} \left\{ 1 - \frac{1}{1 + e^{-x_i^T \theta}} \right\}^{1-y_i} \right\}$$
$$J(\theta) = -\left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

This cost function is called cross-entropy.

Why?

Interpretation of Cross-Entropy Cost Function

Interpretation of Cross-Entropy Cost Function

What is the interpretation of the cost function?

Interpretation of Cross-Entropy Cost Function

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

Interpretation of Cross-Entropy Cost Function

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

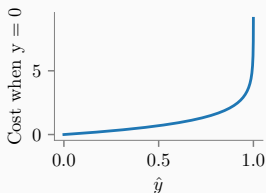
Interpretation of Cross-Entropy Cost Function

What is the interpretation of the cost function?

Let us try to write the cost function for a single example:

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

First, assume y_i is 0, then if \hat{y}_i is 0, the loss is 0; but, if \hat{y}_i is 1, the loss tends towards infinity!



Interpretation of Cross-Entropy Cost Function

Interpretation of Cross-Entropy Cost Function

What is the interpretation of the cost function?

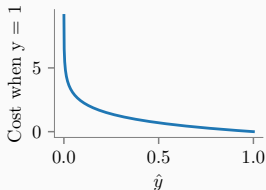
$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

Interpretation of Cross-Entropy Cost Function

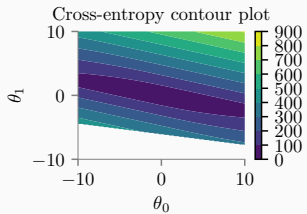
What is the interpretation of the cost function?

$$J(\theta) = -y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)$$

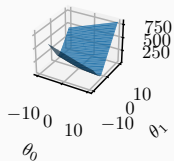
Now, assume y_i is 1, then if \hat{y}_i is 0, the loss is huge; but, if \hat{y}_i is 1, the loss is zero!



Cost function convexity



Cross-entropy surface plot



$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= -\frac{\partial}{\partial \theta_j} \left\{ \sum_{i=1}^n y_i \log(\sigma_\theta(x_i)) + (1 - y_i) \log(1 - \sigma_\theta(x_i)) \right\} \\ &= -\sum_{i=1}^n \left[y_i \frac{\partial}{\partial \theta_j} \log(\sigma_\theta(x_i)) + (1 - y_i) \frac{\partial}{\partial \theta_j} \log(1 - \sigma_\theta(x_i)) \right]\end{aligned}$$

Learning Parameters

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= - \sum_{i=1}^n \left[y_i \frac{\partial}{\partial \theta_j} \log(\sigma_\theta(x_i)) + (1 - y_i) \frac{\partial}{\partial \theta_j} \log(1 - \sigma_\theta(x_i)) \right] \\ &= - \sum_{i=1}^n \left[\frac{y_i}{\sigma_\theta(x_i)} \frac{\partial}{\partial \theta_j} \sigma_\theta(x_i) + \frac{1 - y_i}{1 - \sigma_\theta(x_i)} \frac{\partial}{\partial \theta_j} (1 - \sigma_\theta(x_i)) \right] \quad (1)\end{aligned}$$

Aside:

$$\begin{aligned}\frac{\partial}{\partial z} \sigma(z) &= \frac{\partial}{\partial z} \frac{1}{1 + e^{-z}} = -(1 + e^{-z})^{-2} \frac{\partial}{\partial z} (1 + e^{-z}) \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} = \left(\frac{1}{1 + e^{-z}} \right) \left(\frac{e^{-z}}{1 + e^{-z}} \right) = \sigma(z) \left\{ \frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right\} \\ &= \sigma(z)(1 - \sigma(z))\end{aligned}$$

Learning Parameters

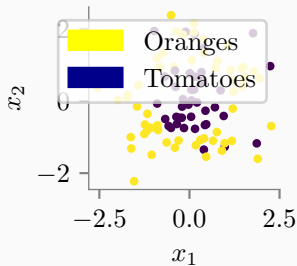
Resuming from (1)

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= - \sum_{i=1}^n \left[\frac{y_i}{\sigma_{\theta}(x_i)} \frac{\partial}{\partial \theta_j} \sigma_{\theta}(x_i) + \frac{1 - y_i}{1 - \sigma_{\theta}(x_i)} \frac{\partial}{\partial \theta_j} (1 - \sigma_{\theta}(x_i)) \right] \\&= - \sum_{i=1}^n \left[\frac{y_i \sigma_{\theta}(x_i)}{\sigma_{\theta}(x_i)} (1 - \sigma_{\theta}(x_i)) \frac{\partial}{\partial \theta_j} (x_i \theta) + \frac{1 - y_i}{1 - \sigma_{\theta}(x_i)} (1 - \sigma_{\theta}(x_i)) \frac{\partial}{\partial \theta_j} (1 - \sigma_{\theta}(x_i)) \right] \\&= - \sum_{i=1}^n \left[y_i (1 - \sigma_{\theta}(x_i)) x_i^j - (1 - y_i) \sigma_{\theta}(x_i) x_i^j \right] \\&= - \sum_{i=1}^n \left[(y_i - y_i \sigma_{\theta}(x_i) - \sigma_{\theta}(x_i) + y_i \sigma_{\theta}(x_i)) x_i^j \right] \\&= \sum_{i=1}^n \left[\sigma_{\theta}(x_i) - y_i \right] x_i^j\end{aligned}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^N [\sigma_{\theta}(x_i) - y_i] x_i^j$$

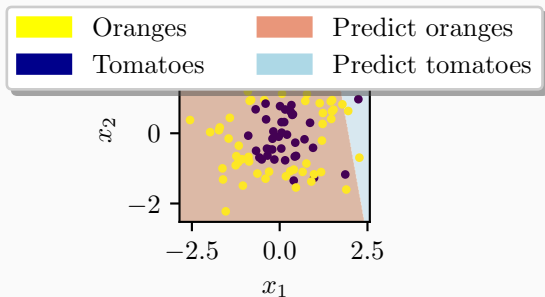
Now, just use Gradient Descent!

Logistic Regression with feature transformation



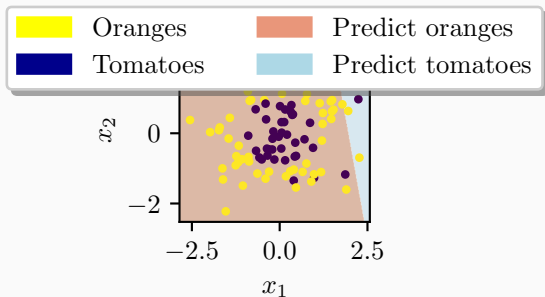
What happens if you apply logistic regression on the above data?

Logistic Regression with feature transformation



Linear boundary will not be accurate here. What is the technical name of the problem?

Logistic Regression with feature transformation

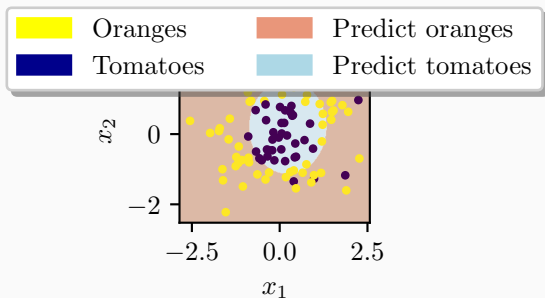


Linear boundary will not be accurate here. What is the technical name of the problem? Bias!

Logistic Regression with feature transformation

$$\phi(x) = \begin{bmatrix} \phi_0(x) \\ \phi_1(x) \\ \vdots \\ \phi_{K-1}(x) \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x^{K-1} \end{bmatrix} \in \mathbb{R}^K$$

Logistic Regression with feature transformation



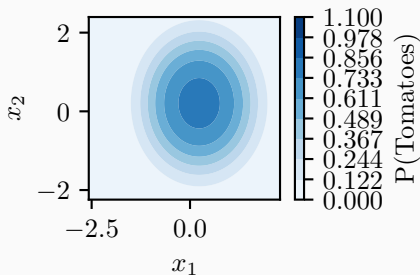
Using x_1^2, x_2^2 as additional features, we are able to learn a more accurate classifier.

Logistic Regression with feature transformation

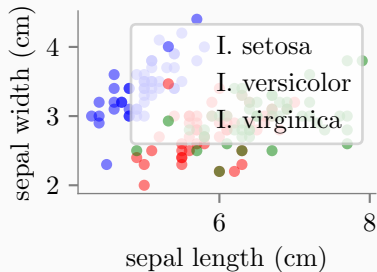
How would you expect the probability contours look like?

Logistic Regression with feature transformation

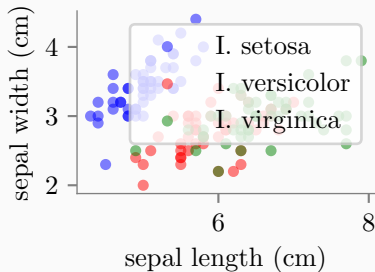
How would you expect the probability contours look like?



Multi-Class Prediction

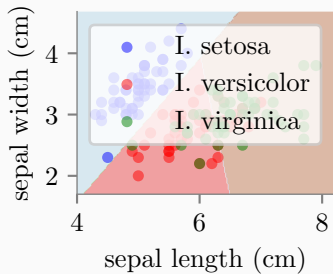


Multi-Class Prediction

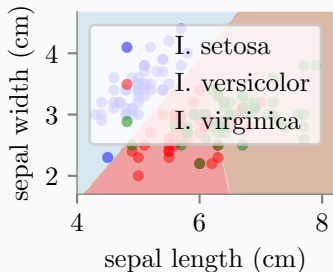


How would you learn a classifier? Or, how would you expect the classifier to learn decision boundaries?

Multi-Class Prediction

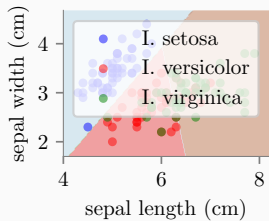


Multi-Class Prediction

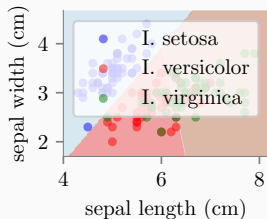


1. Use one-vs.-all on Binary Logistic Regression
2. Use one-vs.-one on Binary Logistic Regression
3. Extend Binary Logistic Regression to Multi-Class Logistic Regression

Multi-Class Prediction

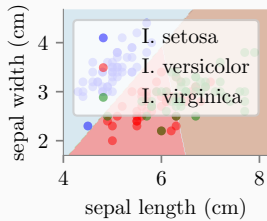


Multi-Class Prediction

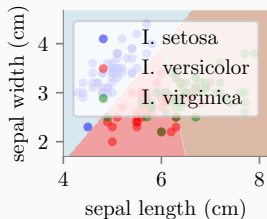


1. Learn $P(\text{setosa (class 1)}) = \mathcal{F}(X\theta_1)$
2. $P(\text{versicolor (class 2)}) = \mathcal{F}(X\theta_2)$
3. $P(\text{virginica (class 3)}) = \mathcal{F}(X\theta_3)$
4. Goal: Learn $\theta_i \forall i \in \{1, 2, 3\}$
5. Question: What could be an \mathcal{F} ?

Multi-Class Prediction



Multi-Class Prediction



1. Question: What could be an \mathcal{F} ?
2. Property: $\sum_{i=1}^3 \mathcal{F}(X\theta_i) = 1$
3. Also $\mathcal{F}(z) \in [0, 1]$
4. Also, $\mathcal{F}(z)$ has squashing properties: $R \mapsto [0, 1]$

$$\begin{aligned} Z &\in \mathbb{R}^d \\ \mathcal{F}(z_i) &= \frac{e^{z_i}}{\sum_{i=1}^d e^{z_i}} \\ \therefore \sum \mathcal{F}(z_i) &= 1 \end{aligned}$$

$\mathcal{F}(z_i)$ refers to probability of class i

Softmax for Multi-Class Logistic Regression

$k = \{1, \dots, K\}$ classes

$$\theta = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \theta_1 & \theta_2 & \cdots & \theta_K \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$$P(y = k|X, \theta) = \frac{e^{X\theta_k}}{\sum_{k=1}^K e^{X\theta_k}}$$

Softmax for Multi-Class Logistic Regression

For $K = 2$ classes,

$$P(y = k|X, \theta) = \frac{e^{X\theta_k}}{\sum_{k=1}^K e^{X\theta_k}}$$

$$P(y = 0|X, \theta) = \frac{e^{X\theta_0}}{e^{X\theta_0} + e^{X\theta_1}}$$

$$\begin{aligned} P(y = 1|X, \theta) &= \frac{e^{X\theta_1}}{e^{X\theta_0} + e^{X\theta_1}} = \frac{e^{X\theta_1}}{e^{X\theta_1} \{1 + e^{X(\theta_0 - \theta_1)}\}} \\ &= \frac{1}{1 + e^{-X\theta'}} \\ &= \text{Sigmoid!} \end{aligned}$$

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

$$= -(0 \times \log(0.1) + 1 \times \log(0.8) + 0 \times \log(0.1))$$

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

$$= -(0 \times \log(0.1) + 1 \times \log(0.8) + 0 \times \log(0.1))$$

Tends to zero

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.3 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.3 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.3 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

$$= -(0 \times \log(0.1) + 1 \times \log(0.4) + 0 \times \log(0.1))$$

Multi-Class Logistic Regression Cost

Assume our prediction and ground truth for the three classes for i^{th} point is:

$$\hat{y}_i = \begin{bmatrix} 0.3 \\ 0.4 \\ 0.3 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

meaning the true class is Class #2

Let us calculate $-\sum_{k=1}^3 y_i^k \log \hat{y}_i^k$

$$= -(0 \times \log(0.1) + 1 \times \log(0.4) + 0 \times \log(0.1))$$

High number! Huge penalty for misclassification!

Multi-Class Logistic Regression Cost

For 2 class we had:

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

Multi-Class Logistic Regression Cost

For 2 class we had:

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

More generally,

Multi-Class Logistic Regression Cost

For 2 class we had:

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

More generally,

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

Multi-Class Logistic Regression Cost

For 2 class we had:

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

More generally,

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

$$J(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right\}$$

Extend to K-class:

$$J(\theta) = - \left\{ \sum_{i=1}^n \sum_{k=1}^K y_i^k \log(\hat{y}_i^k) \right\}$$

Multi-Class Logistic Regression Cost

Now:

$$\frac{\partial J(\theta)}{\partial \theta_k} = \sum_{i=1}^n \left[x_i \left\{ I(y_i = k) - P(y_i = k | x_i, \theta) \right\} \right]$$

Hessian Matrix

The Hessian matrix of $f(\cdot)$ with respect to θ , written $\nabla_{\theta}^2 f(\theta)$ or simply as \mathbb{H} , is the $d \times d$ matrix of partial derivatives,

$$\nabla_{\theta}^2 f(\theta) = \begin{bmatrix} \frac{\partial^2 f(\theta)}{\partial \theta_1^2} & \frac{\partial^2 f(\theta)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 f(\theta)}{\partial \theta_1 \partial \theta_n} \\ \frac{\partial^2 f(\theta)}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 f(\theta)}{\partial \theta_2^2} & \cdots & \frac{\partial^2 f(\theta)}{\partial \theta_2 \partial \theta_n} \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial^2 f(\theta)}{\partial \theta_n \partial \theta_1} & \frac{\partial^2 f(\theta)}{\partial \theta_n \partial \theta_2} & \cdots & \frac{\partial^2 f(\theta)}{\partial \theta_n^2} \end{bmatrix}$$

Newton's Algorithm

The most basic second-order optimization algorithm is Newton's algorithm, which consists of updates of the form,

$$\theta_{k+1} = \theta_k - \mathbb{H}_k^{-1} g_k$$

where g_k is the gradient at step k . This algorithm is derived by making a second-order Taylor series approximation of $f(\theta)$ around θ_k :

$$f_{quad}(\theta) = f(\theta_k) + g_k^T (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T \mathbb{H}_k (\theta - \theta_k)$$

differentiating and equating to zero to solve for θ_{k+1} .

Learning Parameters

Now assume:

$$g(\theta) = \sum_{i=1}^n \left[\sigma_{\theta}(x_i) - y_i \right] x_i^j = X^T (\sigma_{\theta}(X) - y)$$

$$\pi_i = \sigma_{\theta}(x_i)$$

Let \mathbb{H} represent the Hessian of $J(\theta)$

$$\begin{aligned} \mathbb{H} &= \frac{\partial}{\partial \theta} g(\theta) = \frac{\partial}{\partial \theta} \sum_{i=1}^n \left[\sigma_{\theta}(x_i) - y_i \right] x_i^j \\ &= \sum_{i=1}^n \left[\frac{\partial}{\partial \theta} \sigma_{\theta}(x_i) x_i^j - \frac{\partial}{\partial \theta} y_i x_i^j \right] \\ &= \sum_{i=1}^n \sigma_{\theta}(x_i) (1 - \sigma_{\theta}(x_i)) x_i x_i^T \\ &= X^T \text{diag}(\sigma_{\theta}(x_i) (1 - \sigma_{\theta}(x_i))) X \end{aligned}$$

Iteratively reweighted least squares (IRLS)

For binary logistic regression, recall that the gradient and Hessian of the negative log-likelihood are given by:

$$g(\theta)_k = X^T(\pi_k - y)$$

$$H_k = X^T S_k X$$

$$S_k = \text{diag}(\pi_{1k}(1 - \pi_{1k}), \dots, \pi_{nk}(1 - \pi_{nk}))$$

$$\pi_{ik} = \text{sigm}(x_i \theta_k)$$

The Newton update at iteration $k + 1$ for this model is as follows:

$$\begin{aligned}\theta_{k+1} &= \theta_k - H^{-1} g_k \\ &= \theta_k + (X^T S_k X)^{-1} X^T (y - \pi_k) \\ &= (X^T S_k X)^{-1} [(X^T S_k X) \theta_k + X^T (y - \pi_k)] \\ &= (X^T S_k X)^{-1} X^T [S_k X \theta_k + y - \pi_k]\end{aligned}$$

Regularized Logistic Regression

Unregularised:

$$J_1(\theta) = - \left\{ \sum_{i=1}^n y_i \log(\sigma_{\theta}(x_i)) + (1 - y_i) \log(1 - \sigma_{\theta}(x_i)) \right\}$$

L2 Regularization:

$$J(\theta) = J_1(\theta) + \lambda \theta^T \theta$$

L1 Regularization:

$$J(\theta) = J_1(\theta) + \lambda |\theta|$$