


CSC 311: Introduction to Machine Learning

Lecture 5 - Linear Models III

Michael Zhang Chandra Gummaluru

University of Toronto, Winter 2023

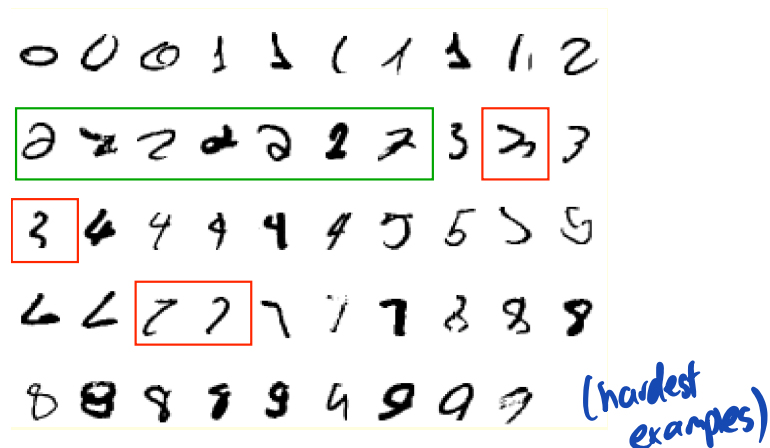
Outline

- 1 Softmax Regression
 - 2 Tracking Model Performance
 - 3 Limits of Linear Classification
 - 4 Midterm Review
 - 5 Introducing Neural Networks
 - 6 Expressivity of a Neural Network
- 

- 1 Softmax Regression
- 2 Tracking Model Performance
- 3 Limits of Linear Classification
- 4 Midterm Review
- 5 Introducing Neural Networks
- 6 Expressivity of a Neural Network

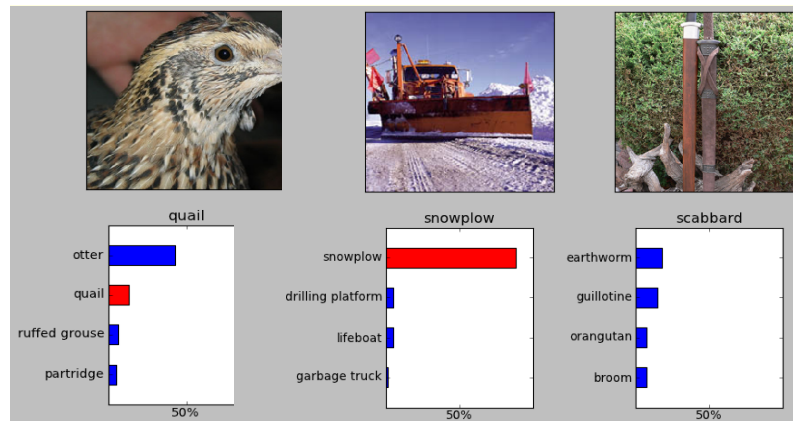
Multi-class Classification

Task is to predict a discrete (> 2)-valued target.



MNIST

FashionMNIST



CIFAR-10

-100

ImageNet (1000 classes)

Targets in Multi-class Classification

- Targets form a discrete set $\{1, \dots, K\}$.
- Represent targets as **one-hot vectors** or **one-of-K encoding**:

$$\mathbf{t} = \underbrace{(0, \dots, 0, 1, 0, \dots, 0)}_{\text{entry } k \text{ is } 1} \in \mathbb{R}^K$$

$N \times D$

Output space

K -dim

$0 \times (K+1)$

\uparrow
bias

Linear Function of Inputs

Vectorized form:

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \text{ or}$$

$$\mathbf{z} = \mathbf{W}\mathbf{x} \text{ with dummy } x_0 = 1$$

Non-vectorized form:

K: number of classes *X.T*
D: number of features

$$z_k = \sum_{j=1}^D w_{kj} x_j + b_k \text{ for } k = 1, 2, \dots, K$$

output x input

- \mathbf{W} : $K \times D$ matrix.
- \mathbf{x} : $D \times 1$ vector.
- \mathbf{b} : $K \times 1$ vector.
- \mathbf{z} : $K \times 1$ vector.

$$\begin{matrix} K \\ \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \end{matrix} \begin{matrix} D \\ \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \end{matrix} + \begin{matrix} K \\ \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \end{matrix} = \begin{matrix} K \\ \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \end{matrix}$$

w *x* *b*

$K \times 1 + K \times 1 = K \times 1$

Generating a Prediction

Interpret z_k as how much the model prefers the k -th prediction.

$$y_i = \begin{cases} 1, & \text{if } i = \arg \max_k z_k \\ 0, & \text{otherwise} \end{cases} \quad \begin{array}{l} \text{(similar to 0-1)} \\ \text{not easily} \\ \text{differentiable} \end{array}$$

How does the $K = 2$ case relate to the binary linear classifiers?

Softmax Regression

- Soften the predictions for optimization.
- A natural activation function is the **softmax function**, a generalization of the logistic function: **np.exp**

$$z = [-2, 0, 3, \dots]$$

$$e^{-2} e^0 e^3, \dots$$

$$y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}}$$

normalize

- Inputs z_k are called the logits.
- Interpret outputs as probabilities. **Y**
- If z_k is much larger than the others, then $\text{softmax}(\mathbf{z})_k \approx 1$ and it behaves like argmax.

$$z_1 = w_1 x$$

$$z_2 = w_2 x$$

$$y_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} = \frac{e^{w_1 x}}{e^{w_1 x} + e^{w_2 x}}$$

What does the $K = 2$ case look like?

Cross-Entropy as Loss Function



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

correct pred
loss is 0
predict 0.99

penalty \gg
predict 0.8

$$\frac{1}{1 + e^{w_2 x - w_1 x}}$$

$$\frac{1}{1 + e^{(w_2 - w_1)x}}$$

Single matrix W in logistic regression

Use cross-entropy as the loss function.

$$\mathcal{L}_{CE}(\mathbf{y}, \mathbf{t}) = - \sum_{k=1}^K t_k \log y_k = -\mathbf{t}^T (\log \mathbf{y}),$$

where the log is applied element-wise.

Often use a combined **softmax-cross-entropy** function.

\hookrightarrow information theory meaning

resume: 12:20

Gradient Descent Updates for Softmax Regression

Softmax Regression:

$$\mathbf{z} = \mathbf{W}\mathbf{x}$$

$$\mathbf{y} = \text{softmax}(\mathbf{z})$$

$$\mathcal{L}_{\text{CE}} = -\mathbf{t}^\top (\log \mathbf{y})$$

Gradient Descent Updates:

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{w}_k} = \frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_k} \cdot \frac{\partial z_k}{\partial \mathbf{w}_k} = (y_k - t_k) \cdot \mathbf{x}$$

$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \alpha \frac{1}{N} \sum_{i=1}^N (y_k^{(i)} - t_k^{(i)}) \mathbf{x}^{(i)}$$

$$(X^T X + \lambda I) w^* = X^T t \quad \leftarrow \text{direct solution}$$

↑
matrix
invertible

problem $\|Xw - t\|_2^2 + \lambda \|w\|_2^2$

$$\|w\|_2^2 = w_1^2 + w_2^2 + \dots + w_d^2 \quad (\text{linear eq.})$$

$$= w^T w = w \cdot w$$

$$\nabla_w \|w\|_2^2 = 2w$$

$$f(x) = X^T A x$$

quadratic form

Multivariate Gaussians

$$w^T w = w^T I w$$

$$\nabla_w w^T I w = 2I w = 2w$$

$$\mathbb{R}^d \rightarrow \mathbb{R}$$

$$\begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \vdots \end{bmatrix} = \begin{bmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_d \end{bmatrix}$$

$$= 2w$$

$$\propto e^{-\frac{x^2}{2}}$$

$$x: \mathbb{R}^n$$

f ↓

\mathbb{R}

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$= a_{11} x_1^2 + a_{12} x_1 x_2$$

$$+ a_{21} x_2 x_1 + a_{22} x_2^2$$

$$= \sum_i \sum_j A_{ij} x_i x_j$$

x_2

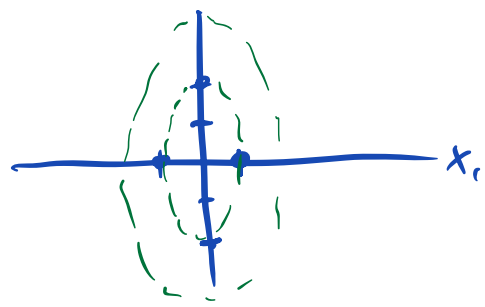
$$x^T \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} x$$

$$= 4x_1^2 + x_2^2 = r$$

$$r = 4$$

$$x_1 = \pm 1$$

$$x_2 = \pm 2$$



$$\begin{bmatrix} x \\ \vdots \\ x_i \\ \vdots \\ x \end{bmatrix}^T \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & & & \\ \vdots & & & \\ a_{ij} & & & \\ \vdots & & & \\ a_{in} & & & \\ a_{nn} & & & \end{bmatrix} \begin{bmatrix} x \\ \vdots \\ x \\ \vdots \\ x \end{bmatrix}$$

$$\frac{\partial f}{\partial x_i} = d \left[\sum_{j \neq i} a_{ij} x_i x_j + \sum_{j \neq i} a_{ji} x_i x_j + a_{ii} x_i^2 \right]$$

dx_i

$2Ax$ symmetric
 $(A+A^T)x$ not sym.

$$= \sum_{j \neq i} a_{ij} x_j + \sum_{j \neq i} a_{ji} x_j + 2a_{ii} x_i$$

$$= \sum_{j=1}^n a_{ij} x_j + \sum_{j=1}^n a_{ji} x_j$$