Multi-class linear prediction

COMS 4771 Fall 2025

Multi-class prediction

Multi-class prediction usually means classification problems with $|\mathcal{Y}| \geq 3$ (... but technically also includes $|\mathcal{Y}| = 2$)

- ► Assume we use zero-one loss, so corresponding risk is error rate
- ightharpoonup Best prediction of Y given X=x is

$$f^{\star}(x) = \underset{k \in \mathcal{Y}}{\operatorname{arg\,max}} \Pr(Y = k \mid X = x)$$

	$ \mathcal{Y} = 2$	$ \mathcal{Y} \ge 3$
nearest neighbors	\checkmark	\checkmark
decision trees	\checkmark	\checkmark
generative models	√	√
logistic regression	√	?
Perceptron	√	?

Linear classifiers are inherently binary output functions; need some work to extend to $|\mathcal{Y}| \geq 3$

1/9

Multi-class logistic regression

Generalize logistic regression model to K classes, $[K] = \{1, 2, \dots, K\}$

$$\Pr(Y = k \mid X = x) = \frac{\exp(x^{\mathsf{T}} w^{(k)})}{\sum_{c \in [K]} \exp(x^{\mathsf{T}} w^{(c)})}$$

 $\blacktriangleright \ K$ weight vectors $w^{(1)}, \dots, w^{(K)} \in \mathbb{R}^d$ are parameters of the model

$$w^{(k)} = (w_1^{(k)}, \dots, w_d^{(k)})$$

Total of Kd parameters

In this model with parameters $w^{(1)},\dots,w^{(K)}$, classifier with smallest error rate is

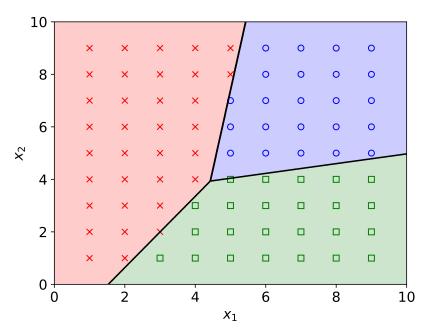
$$f^{\star}(x) = \operatorname*{arg\,max}_{k \in [K]} \Pr(Y = k \mid X = x) = \operatorname*{arg\,max}_{k \in [K]} x^{\mathsf{\scriptscriptstyle T}} w^{(k)}$$

Maximum likelihood estimation: given training data S from $\mathbb{R}^d \times [K]$, log-likelihood of $w^{(1)},\dots,w^{(K)}$ is

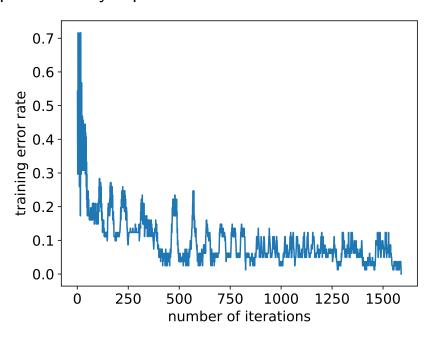
$$\ln L(w^{(1)}, \dots, w^{(K)}) = \sum_{(x,y)\in\mathbb{S}} \ln \left(\frac{\exp(x^{\mathsf{T}} w^{(y)})}{\sum_{c\in[K]} \exp(x^{\mathsf{T}} w^{(c)})} \right)$$
$$= \sum_{(x,y)\in\mathbb{S}} \left(x^{\mathsf{T}} w^{(y)} - \ln \left(\sum_{c\in[K]} \exp(x^{\mathsf{T}} w^{(c)}) \right) \right)$$

▶ Negative log-likelihood turns out to be a convex objective function

Synthetic example: "linearly separable" dataset



Synthetic example: "linearly separable" dataset



6/9

Can also interpret negative log-likelihood as sum of log losses, like in binary case, for prediction function

$$\hat{p}_W(x) = \operatorname{softmax}(Wx)$$

where $\operatorname{softmax} \colon \mathbb{R}^K \to \mathbb{R}^K$ (soft(arg)max function) is defined by

$$\operatorname{softmax}(u)_k = \frac{\exp(u_k)}{\sum_{c \in [K]} \exp(u_c)}$$

- $lackbox{W} \in \mathbb{R}^{K imes d}$ is matrix of parameters, one row per class
- ► Negative log-likelihood is

$$J(W) = \sum_{(x,y)\in\mathcal{S}} -\ln(\hat{p}_W(x)_y)$$

Setup

Gradient descent code

```
for t in range(num_iter):
J = loss(f(x), y)
J.backward()
with torch.no_grad():
    W -= eta * W.grad
    W.grad.zero_()
```

8/9

Example: iris dataset

▶ Multi-class logistic regression ($\eta = 0.01$, $T = 2^{17}$ iterations)

► Test error rate: 3.33%

