

CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

Feb 5, 2021

1 / 31

Logistics

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Multiclass Classification

3 / 31

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Multiclass Classification

2 / 31

Logistics

Logistics

- HW 1 is due today, and HW 2 will be assigned.
- Please form the groups for the project, we'll have groups of 3 students working together. Use piazza to find group members.

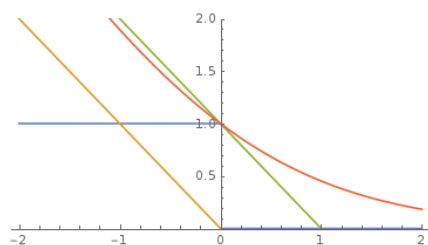
4 / 31

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Multiclass Classification

5 / 31

Step 2. Pick the **surrogate loss**



- **perceptron loss** $\ell_{\text{perceptron}}(z) = \max\{0, -z\}$ (used in Perceptron)
- **hinge loss** $\ell_{\text{hinge}}(z) = \max\{0, 1 - z\}$ (used in SVM and many others)
- **logistic loss** $\ell_{\text{logistic}}(z) = \log(1 + \exp(-z))$ (used in logistic regression)

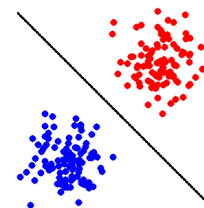
7 / 31

Summary

Linear models for **binary** classification:

Step 1. Model is the set of **separating hyperplanes**

$$\mathcal{F} = \{f(x) = \text{sgn}(w^T x) \mid w \in \mathbb{R}^D\}$$



6 / 31

Step 3. Find empirical risk minimizer (ERM):

$$w^* = \underset{w \in \mathbb{R}^D}{\text{argmin}} F(w) = \underset{w \in \mathbb{R}^D}{\text{argmin}} \frac{1}{N} \sum_{n=1}^N \ell(y_n w^T x_n)$$

using

- **GD:** $w \leftarrow w - \eta \nabla F(w)$
- **SGD:** $w \leftarrow w - \eta \tilde{\nabla} F(w)$
- **Newton:** $w \leftarrow w - (\nabla^2 F(w))^{-1} \nabla F(w)$

8 / 31

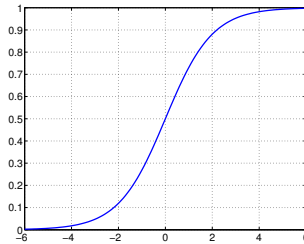
A Probabilistic view of logistic regression

Minimizing logistic loss = MLE for the sigmoid model

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w})$$

where

$$\mathbb{P}(y | \mathbf{x}; \mathbf{w}) = \sigma(y \mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-y \mathbf{w}^T \mathbf{x}}}$$



9 / 31

Outline

- 1 Logistics
- 2 Review of Last Lecture
- 3 Multiclass Classification
 - Multinomial logistic regression
 - Reduction to binary classification

10 / 31

Classification

Recall the setup:

- input (feature vector): $\mathbf{x} \in \mathbb{R}^D$
- output (label): $y \in [C] = \{1, 2, \dots, C\}$
- goal: learn a mapping $f: \mathbb{R}^D \rightarrow [C]$

Examples:

- recognizing digits ($C = 10$) or letters ($C = 26$ or 52)
- predicting weather: sunny, cloudy, rainy, etc
- predicting image category: ImageNet dataset ($C \approx 20K$)

Nearest Neighbor Classifier naturally works for arbitrary C .

11 / 31

Linear models: from binary to multiclass

Step 1: *What should a linear model look like for multiclass tasks?*

Note: a linear model for binary tasks (switching from $\{-1, +1\}$ to $\{1, 2\}$)

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 2 & \text{if } \mathbf{w}^T \mathbf{x} < 0 \end{cases}$$

can be written as

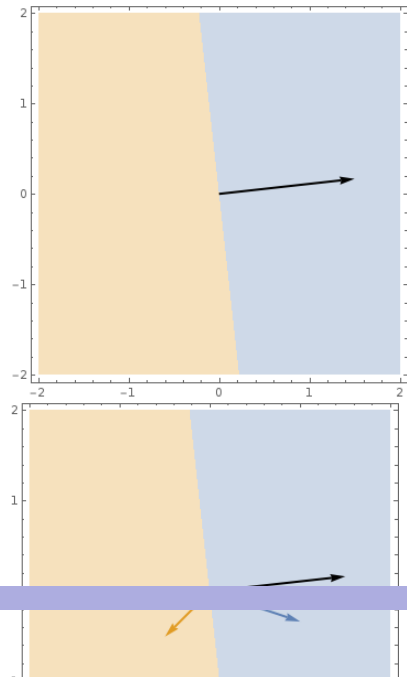
$$\begin{aligned} f(\mathbf{x}) &= \begin{cases} 1 & \text{if } \mathbf{w}_1^T \mathbf{x} \geq \mathbf{w}_2^T \mathbf{x} \\ 2 & \text{if } \mathbf{w}_2^T \mathbf{x} > \mathbf{w}_1^T \mathbf{x} \end{cases} \\ &= \underset{k \in \{1, 2\}}{\operatorname{argmax}} \mathbf{w}_k^T \mathbf{x} \end{aligned}$$

for any $\mathbf{w}_1, \mathbf{w}_2$ s.t. $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$

Think of $\mathbf{w}_k^T \mathbf{x}$ as **a score for class k** .

12 / 31

Linear models: from binary to multiclass



$$\mathbf{w} = \left(\frac{3}{2}, \frac{1}{6}\right) = \mathbf{w}_1 - \mathbf{w}_2$$

$$\mathbf{w}_1 = \left(1, -\frac{1}{3}\right)$$

$$\mathbf{w}_2 = \left(-\frac{1}{2}, -\frac{1}{2}\right)$$

- Blue class:

$$\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$$

Linear models for multiclass classification

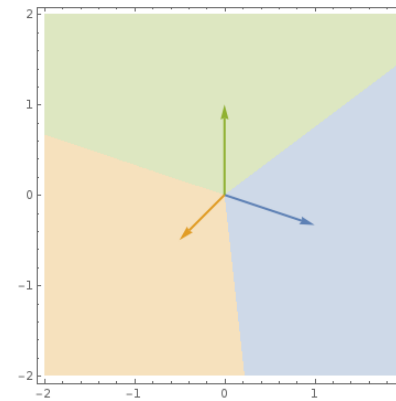
$$\mathcal{F} = \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} \mathbf{w}_k^T \mathbf{x} \mid \mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^D \right\}$$

$$= \left\{ f(\mathbf{x}) = \operatorname{argmax}_{k \in [C]} (\mathbf{W}\mathbf{x})_k \mid \mathbf{W} \in \mathbb{R}^{C \times D} \right\}$$

Step 2: *How do we generalize perceptron/hinge/logistic loss?*

This lecture: focus on the more popular **logistic loss**

Linear models: from binary to multiclass



$$\mathbf{w}_1 = \left(1, -\frac{1}{3}\right)$$

$$\mathbf{w}_2 = \left(-\frac{1}{2}, -\frac{1}{2}\right)$$

$$\mathbf{w}_3 = (0, 1)$$

- Blue class:
 $\{\mathbf{x} : 1 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Orange class:
 $\{\mathbf{x} : 2 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$
- Green class:
 $\{\mathbf{x} : 3 = \operatorname{argmax}_k \mathbf{w}_k^T \mathbf{x}\}$

Multinomial logistic regression: a probabilistic view

Observe: for binary logistic regression, with $\mathbf{w} = \mathbf{w}_1 - \mathbf{w}_2$:

$$\mathbb{P}(y = 1 \mid \mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \frac{e^{\mathbf{w}_1^T \mathbf{x}}}{e^{\mathbf{w}_1^T \mathbf{x}} + e^{\mathbf{w}_2^T \mathbf{x}}} \propto e^{\mathbf{w}_1^T \mathbf{x}}$$

Naturally, for multiclass:

$$\mathbb{P}(y = k \mid \mathbf{x}; \mathbf{W}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}}}{\sum_{k' \in [C]} e^{\mathbf{w}_{k'}^T \mathbf{x}}} \propto e^{\mathbf{w}_k^T \mathbf{x}}$$

This is called the *softmax function*.

Applying MLE again

Maximize probability of seeing labels y_1, \dots, y_N given $\mathbf{x}_1, \dots, \mathbf{x}_N$

$$P(\mathbf{W}) = \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{W}) = \prod_{n=1}^N \frac{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}}{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\mathbf{W}) = \sum_{n=1}^N \ln \left(\frac{\sum_{k \in [C]} e^{\mathbf{w}_k^T \mathbf{x}_n}}{e^{\mathbf{w}_{y_n}^T \mathbf{x}_n}} \right) = \sum_{n=1}^N \ln \left(1 + \sum_{k \neq y_n} e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)$$

This is the **multiclass logistic loss**, a.k.a **cross-entropy loss**.

When $C = 2$, this is the same as binary logistic loss.

17 / 31

Step 3: Optimization

Apply **SGD**: what is the gradient of

$$g(\mathbf{W}) = \ln \left(1 + \sum_{k' \neq y_n} e^{(\mathbf{w}_{k'} - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right) ?$$

18 / 31

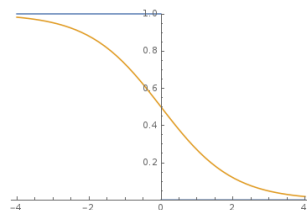
SGD for Binary Classification case (last lecture)

Recall that $\ell_{\text{logistic}}(z) = \ln(1 + \exp(-z))$

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \eta \tilde{\nabla} F(\mathbf{w}) \\ &= \mathbf{w} - \eta \nabla_{\mathbf{w}} \ell_{\text{logistic}}(y_n \mathbf{w}^T \mathbf{x}_n) \quad (n \in [N] \text{ is drawn u.a.r.}) \\ &= \mathbf{w} - \eta \left(\frac{\partial \ell_{\text{logistic}}(z)}{\partial z} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\ &= \mathbf{w} - \eta \left(\frac{-e^{-z}}{1 + e^{-z}} \Big|_{z=y_n \mathbf{w}^T \mathbf{x}_n} \right) y_n \mathbf{x}_n \\ &= \mathbf{w} + \eta \sigma(-y_n \mathbf{w}^T \mathbf{x}_n) y_n \mathbf{x}_n \\ &= \mathbf{w} + \eta \mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w}) y_n \mathbf{x}_n \end{aligned}$$

This is a **soft version of Perceptron!**

$\mathbb{P}(-y_n | \mathbf{x}_n; \mathbf{w})$ versus $\mathbb{I}[y_n \neq \text{sgn}(\mathbf{w}^T \mathbf{x}_n)]$



19 / 31

Step 3: Optimization

Apply **SGD**: what is the gradient of

$$g(\mathbf{W}) = \ln \left(1 + \sum_{k' \neq y_n} e^{(\mathbf{w}_{k'} - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right) ?$$

It's a $C \times D$ matrix. Let's focus on the k -th row:

If $k \neq y_n$:

$$\nabla_{\mathbf{w}_k} g(\mathbf{W}) = \frac{e^{(\mathbf{w}_k - \mathbf{w}_{y_n})^T \mathbf{x}_n}}{1 + \sum_{k' \neq y_n} e^{(\mathbf{w}_{k'} - \mathbf{w}_{y_n})^T \mathbf{x}_n}} \mathbf{x}_n^T = \mathbb{P}(k | \mathbf{x}_n; \mathbf{W}) \mathbf{x}_n^T$$

else:

$$\nabla_{\mathbf{w}_k} g(\mathbf{W}) = \frac{-\left(\sum_{k' \neq y_n} e^{(\mathbf{w}_{k'} - \mathbf{w}_{y_n})^T \mathbf{x}_n} \right)}{1 + \sum_{k' \neq y_n} e^{(\mathbf{w}_{k'} - \mathbf{w}_{y_n})^T \mathbf{x}_n}} \mathbf{x}_n^T = (\mathbb{P}(y_n | \mathbf{x}_n; \mathbf{W}) - 1) \mathbf{x}_n^T$$

20 / 31

SGD for multinomial logistic regression

Initialize $\mathbf{W} = \mathbf{0}$ (or randomly). Repeat:

- 1 pick $n \in [N]$ uniformly at random
- 2 update the parameters

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \begin{pmatrix} \mathbb{P}(y = 1 \mid \mathbf{x}_n; \mathbf{W}) \\ \vdots \\ \mathbb{P}(y = y_n \mid \mathbf{x}_n; \mathbf{W}) - 1 \\ \vdots \\ \mathbb{P}(y = C \mid \mathbf{x}_n; \mathbf{W}) \end{pmatrix} \mathbf{x}_n^T$$

Think about why the algorithm makes sense intuitively.

21 / 31

Reduce multiclass to binary

Is there an *even more general and simpler approach* to derive multiclass classification algorithms?

Given a binary classification algorithm (*any one*, not just linear methods), can we turn it to a multiclass algorithm, *in a black-box manner*?

Yes, there are in fact many ways to do it.

- **one-versus-all** (one-versus-rest, one-against-all, etc)
- **one-versus-one** (all-versus-all, etc)
- **Error-Correcting Output Codes** (ECOC)
- **tree-based reduction**

23 / 31

A note on prediction

Having learned \mathbf{W} , we can either

- make a *deterministic* prediction $\arg\max_{k \in [C]} \mathbf{w}_k^T \mathbf{x}$
- make a *randomized* prediction according to $\mathbb{P}(k \mid \mathbf{x}; \mathbf{W}) \propto e^{\mathbf{w}_k^T \mathbf{x}}$

In either case, **(expected) mistake is bounded by logistic loss**

- deterministic

$$\mathbb{I}[f(\mathbf{x}) \neq y] \leq \ln \left(1 + \sum_{k \neq y} e^{(\mathbf{w}_k - \mathbf{w}_y)^T \mathbf{x}} \right)$$

- randomized

$$\mathbb{E} [\mathbb{I}[f(\mathbf{x}) \neq y]] = 1 - \mathbb{P}(y \mid \mathbf{x}; \mathbf{W}) \leq -\ln \mathbb{P}(y \mid \mathbf{x}; \mathbf{W})$$

22 / 31

One-versus-all (OvA)

(picture credit: link)

Idea: train C binary classifiers to learn “**is class k or not?**” for each k .

Training: for each class $k \in [C]$,

- re-label examples with class k as $+1$, and all others as -1
- train a binary classifier h_k using this new dataset

		■	■	■	■
x_1	■	x_1 —	x_1 +	x_1 —	x_1 —
x_2	■	x_2 —	x_2 —	x_2 +	x_2 —
x_3	■	x_3 —	x_3 —	x_3 —	x_3 +
x_4	■	x_4 —	x_4 +	x_4 —	x_4 —
x_5	■	x_5 +	x_5 —	x_5 —	x_5 —
		⇓	⇓	⇓	⇓
		h_1	h_2	h_3	h_4

24 / 31

One-versus-all (OvA)

Prediction: for a new example x

- ask each h_k : **does this belong to class k ?** (i.e. $h_k(x)$)
- randomly pick among all k 's s.t. $h_k(x) = +1$.

Issue: will (probably) make a mistake *as long as one of h_k errs*.

25 / 31

One-versus-one (OvO)

(picture credit: link)

Idea: train $\binom{C}{2}$ binary classifiers to learn “**is class k or k' ?**”.

Training: for each pair (k, k') ,

- re-label class k examples as $+1$ and class k' examples as -1
- discard all other examples
- train a binary classifier $h_{(k,k')}$ using this new dataset

		■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■
		■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■	■ vs. ■
x_1	■	x_1 —		x_2 —	x_1 —		x_1 —
x_2	■		x_2 —	x_2 +		x_3 —	x_2 +
x_3	■			x_3 —	x_3 +	x_3 —	
x_4	■	x_4 —			x_4 —		x_4 —
x_5	■	x_5 +	x_5 +			x_5 +	
		↓	↓	↓	↓	↓	↓
		$h_{(1,2)}$	$h_{(1,3)}$	$h_{(3,4)}$	$h_{(4,2)}$	$h_{(1,4)}$	$h_{(3,2)}$

26 / 31

One-versus-one (OvO)

Prediction: for a new example x

- ask each classifier $h_{(k,k')}$ to **vote for either class k or k'**
- predict the class with the most votes (break tie in some way)

More robust than one-versus-all, but *slower* in prediction.

27 / 31

Error-correcting output codes (ECOC)

(picture credit: link)

Idea: based on a code $M \in \{-1, +1\}^{C \times L}$, train L binary classifiers to learn “**is bit b on or off**”.

Training: for each bit $b \in [L]$

- re-label example x_n as $M_{y_n, b}$
- train a binary classifier h_b using this new dataset.

M	1	2	3	4	5
■	+	—	+	—	+
■	—	—	+	+	+
■	+	+	—	—	—
■	+	+	+	+	—

	1	2	3	4	5
x_1	x_1 —	x_1 —	x_1 +	x_1 +	x_1 +
x_2	x_2 +	x_2 +	x_2 —	x_2 —	x_2 —
x_3	x_3 +	x_3 +	x_3 +	x_3 +	x_3 —
x_4	x_4 —	x_4 —	x_4 +	x_4 +	x_4 +
x_5	x_5 +	x_5 —	x_5 +	x_5 —	x_5 +
	↓	↓	↓	↓	↓
	h_1	h_2	h_3	h_4	h_5

28 / 31

Error-correcting output codes (ECOC)

Prediction: for a new example x

- compute the **predicted code** $c = (h_1(x), \dots, h_L(x))^T$
- predict the class with the **most similar code**: $k = \operatorname{argmax}_k (Mc)_k$

How to design the code M ?

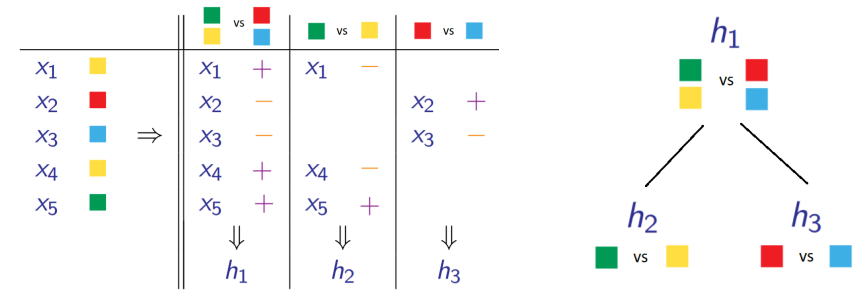
- the more *dissimilar* the codes, the more robust
 - if any two codes are d bits away, then prediction can tolerate about $d/2$ errors
- random code* is often a good choice

29 / 31

Tree based method

Idea: train $\approx C$ binary classifiers to learn “**belongs to which half?**”.

Training: see pictures



Prediction is also natural, *but is very fast!* (think ImageNet where $C \approx 20K$)

30 / 31

Comparisons

In big O notation,

Reduction	#training points	test time	remark
OvA	CN	C	not robust
OvO	CN	C^2	can achieve very small training error
ECOC	LN	L	need diversity when designing code
Tree	$(\log_2 C)N$	$\log_2 C$	good for “extreme classification”

31 / 31