

# Multinomial Regression (MNR)

Tuesday, October 27, 2020 5:05 PM

## Multi-Nomial Regression (MNR) : Intro

- What we have seen so far:

GDA/LDA classifiers: probabilistic classifiers

Based on the assumption that  $\underline{x}|y \sim \mathcal{N}(\mu_y, \Sigma_y)$

These classifiers are Probabilistic Generative Models

- MNR belongs to the Probabilistic Discriminative Model Family

It makes no attempt to model  $\underline{x}$

For GDA/LDA we modeled  $P(\underline{x}, y) = P(\underline{x}|y) P(y)$

$y \sim \text{Cat}(\underline{P})$   
↑  
vector of class prior probabilities

$$= \prod_{k=1}^K P(\underline{x}|\underline{\mu}_k, \Sigma_k)$$

Parenthesis: Categorical Distribution.

$$y \sim \text{Cat}(\underline{P}) \stackrel{\text{def}}{\iff} P\{y=k\} = \begin{cases} p_k & k=1, \dots, K \\ 0 & \text{otherwise} \end{cases}$$

$$\Leftrightarrow P\{y=k\} = \prod_{t=1}^d p_t^{I[t=k]}$$

Finally,  $P(y|\underline{x}) = \frac{P(\underline{x}|y) P(y)}{P(\underline{x})}$

↳ distn. of the mixture.

For MNR, if models  $P(y|\underline{x})$  directly

$$P(y=k|\underline{x}) = \frac{e^{h_k(\underline{x})}}{\sum_{k'=1}^K e^{h_{k'}(\underline{x})}}$$

e.g.  $K=2$   $P(y=1|\underline{x}) = \frac{e^{h_1(\underline{x})}}{e^{h_1(\underline{x})} + e^{h_2(\underline{x})}}$

$$P(y=2|\underline{x}) = \frac{e^{h_2(\underline{x})}}{e^{h_1(\underline{x})} + e^{h_2(\underline{x})}} =$$

$h_k$ : scoring/score function of class  $k$ .

$$P(y=2|x) = \frac{e^{h_2(x)}}{e^{h_1(x)} + e^{h_2(x)}} = \\ = 1 - P(y=1|x)$$

$h_k$ : scoring/score function of class  $k$ .

Interpretation: the higher the value of  $h_k(x)$  is, the more the evidence/likelihood that  $x$  was drawn from class  $k$

$$\begin{bmatrix} P(y=1|x) \\ P(y=2|x) \\ \vdots \\ P(y=G|x) \end{bmatrix} = \text{softmax}(h(x)) = \begin{bmatrix} \text{softmax}_1(h(x)) \\ \vdots \\ \text{softmax}_G(h(x)) \end{bmatrix}$$

$$\text{where } \text{softmax}_k(h(x)) \triangleq \frac{e^{h_k(x)}}{\sum_{k'} e^{h_{k'}(x)}}$$

MNR would predict  $\hat{y} = \arg \max \{P(y=k|x)\}$

Furthermore, MNR use the following form of score functions:

$$h_k(x) = \underline{w}_k^T x \quad k=1, \dots, G$$

MLPs used as classifiers:  $h_k(x)$  is a much more complex, non-linear function of  $x$ , with many more model parameters

Example: MNR for  $G=3$

$$f_1(x) = \frac{e^{\underline{w}_1^T x}}{e^{\underline{w}_1^T x} + e^{\underline{w}_2^T x} + e^{\underline{w}_3^T x}}, \text{ which models/estimates } P(y=1|x)$$

$$f_2(x) = \frac{e^{\underline{w}_2^T x}}{e^{\underline{w}_1^T x} + e^{\underline{w}_2^T x} + e^{\underline{w}_3^T x}} \quad -\text{II-} \quad P(y=2|x)$$

$$f_3(x) = \frac{e^{\underline{w}_3^T x}}{e^{\underline{w}_1^T x} + e^{\underline{w}_2^T x} + e^{\underline{w}_3^T x}} \quad -\text{III-} \quad P(y=3|x)$$

$f_k$ : model's (softmax's)  $k^{\text{th}}$  output.

Model parameters:  $\underline{w}_1, \underline{w}_2, \underline{w}_3$

one can organize these in a single matrix

$$W = \begin{bmatrix} \underline{w}_1^T \\ \underline{w}_2^T \\ \underline{w}_3^T \end{bmatrix} \in \mathbb{R}^{C \times D}$$

$$W = \begin{bmatrix} w_1^T \\ w_2^T \\ w_3^T \end{bmatrix} \in \mathbb{R}^{C \times D}$$

$$\underline{h}(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ h_3(x) \end{bmatrix} = \begin{bmatrix} w_1^T x \\ w_2^T x \\ w_3^T x \end{bmatrix} = \begin{bmatrix} w_1^T \\ w_2^T \\ w_3^T \end{bmatrix} x = Wx$$

- Why softmax(·)

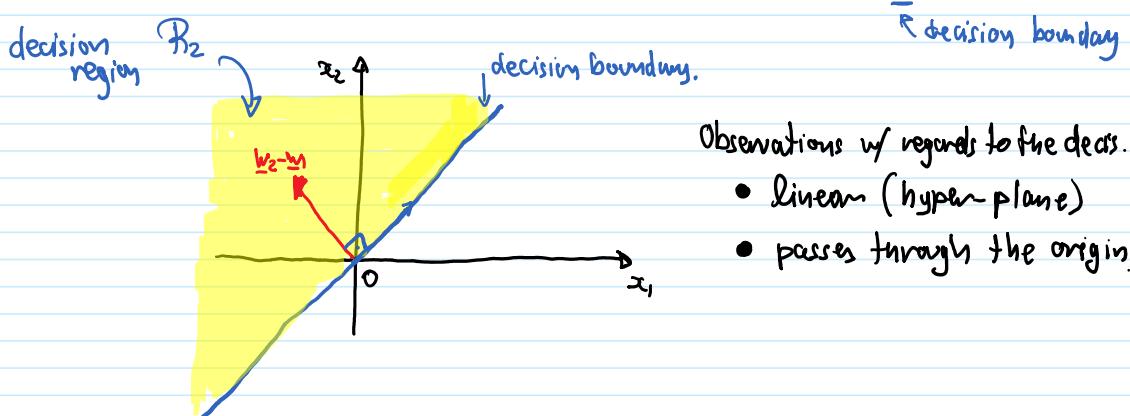
- 1)  $\text{softmax}_k(h(x)) \in (0, 1)$   
 2)  $\sum_k \text{softmax}_k(h(x)) = 1$
- }  $\text{softmax}(·)$  is a good function to model a set of probability of mutually exclusive event

- Decision Regions/Boundaries (for simplicity,  $C=2$ )

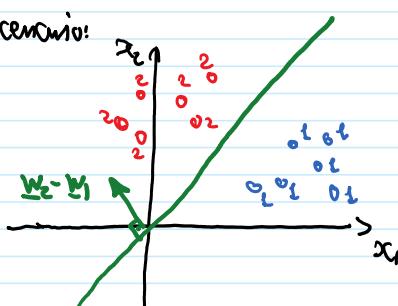
$$\hat{y}=2 \Leftrightarrow s'_2(h(x)) > s'_1(h(x)) \Leftrightarrow$$

$$\Leftrightarrow e^{h_2(x)} > e^{h_1(x)} \Leftrightarrow$$

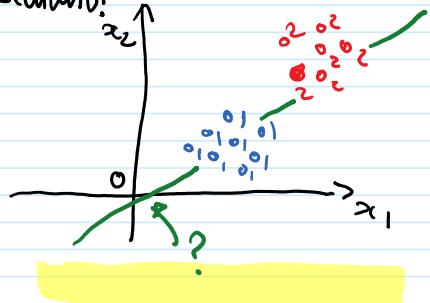
$$\Leftrightarrow h_2(x) > h_1(x) \Leftrightarrow \dots \Leftrightarrow \underbrace{(w_2 - w_1)^T x}_{=0} > 0$$



A good scenario:



A bad scenario:



The work-around: make decision boundaries to be hyperplanes that do not have to pass through the origin.

- "Trick": 1) augment each input  $\underline{x}$  with an extra "dummy" feature 1  
 2) augment each weight vector with an extra bias/intcept parameter

E.g. Use as input  $\tilde{\underline{x}} \triangleq \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} \in \mathbb{R}^{D+1}$  instead of just  $\underline{x} \in \mathbb{R}^D$ .

use weight vectors  $\tilde{\underline{w}}_k \triangleq \begin{bmatrix} \underline{w}_k \\ b_k \end{bmatrix} \in \mathbb{R}^{D+1}$  instead of just  $\underline{w}_k$

Here's why it works (e.g.  $G=2$ )

$$\hat{y}=2 \Leftrightarrow (\tilde{\underline{w}}_2 - \tilde{\underline{w}}_1)^T \tilde{\underline{x}} > 0 \Leftrightarrow \begin{bmatrix} \underline{w}_2 - \underline{w}_1 \\ b_2 - b_1 \end{bmatrix}^T \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix} > 0 \Leftrightarrow$$

$$\Leftrightarrow (\underline{w}_2 - \underline{w}_1)^T \underline{x} + (b_2 - b_1) = 0$$

equation of a hyperplane in general position.

### ► Translation Invariance of decision rule.

From previous example,

$$\hat{y}=2 \Leftrightarrow (\tilde{\underline{w}}_2 - \tilde{\underline{w}}_1)^T \tilde{\underline{x}} > 0$$

Consider the weight translations:  $\tilde{\underline{w}}'_1 \triangleq \tilde{\underline{w}}_1 + \tilde{\underline{w}}_0$  ← arbitrary vector  
 $\tilde{\underline{w}}'_2 \triangleq \tilde{\underline{w}}_2 + \tilde{\underline{w}}_0$

$$(\tilde{\underline{w}}'_2 - \tilde{\underline{w}}'_1)^T \tilde{\underline{x}} > 0 \Leftrightarrow (\tilde{\underline{w}}_2 - \tilde{\underline{w}}_1)^T \tilde{\underline{x}} > 0 \Leftrightarrow \hat{y}=2.$$

In general ( $G \geq 2$ ),  $\begin{bmatrix} \tilde{\underline{w}}_1^T + \tilde{\underline{w}}_0^T \\ \vdots \\ \tilde{\underline{w}}_G^T + \tilde{\underline{w}}_0^T \end{bmatrix}$  and  $\begin{bmatrix} \tilde{\underline{w}}_1^T \\ \vdots \\ \tilde{\underline{w}}_G^T \end{bmatrix}$  produce the same

same decision regions and, hence, label predictions.

## ④ Training MNR models

Goal: to train MNR classifier whose empirical off loss (sample misclassif. rate) is minimum.

$$E_{\text{off}} = \frac{1}{N} \sum_{n=1}^N \ell_{\text{off}}(\hat{y}_n, y_n) \quad ①$$

$$\ell_{\text{off}}(\hat{y}(x), y) \triangleq [\max_{k \neq y} h_k(x) < 0] = \begin{cases} 0 & \text{if } x \text{ was correctly classified} \\ 1 & \text{if } x \text{ was misclassified} \end{cases} \quad ②$$

$$\begin{aligned} \hat{y} &= \arg \max_k \hat{P}(l=k|x) = \arg \max_k \text{softmax}_k(h_k(x)) = \\ &= \arg \max_k e^{h_k(x)} = \arg \max_k h_k(x) \end{aligned}$$

If  $h_y(x)$  is the largest of the  $h_k(x)$ 's: prediction is correct.  
 -1- is not      -1- : wrong prediction.

$$①② \Rightarrow E_{\text{off}} = \frac{1}{N} \sum_{n=1}^N [\max_{k \neq y_n} h_k(x_n) < 0]$$

$\uparrow$  discontinuous  $\Rightarrow$  non-differentiable.

## ► Cross-Entropy (CE) loss & Average CE (ACE)

Instead of the off loss, we are going to use this surrogate loss

$$\ell_{\text{CE}}(t, \hat{s}(x)) = - \sum_{k=1}^C t_k \ln(\hat{s}_k(x)) = - \ln \hat{s}_y(x)$$

$\uparrow$  softmax function evaluated using  $x$   
 1-hot encoding of  $x$ 's label  $y$

$$\hat{s}(x) \triangleq \frac{1}{\sum_{k=1}^C e^{h_k(x)}} \begin{bmatrix} e^{h_1(x)} \\ e^{h_2(x)} \\ \vdots \\ e^{h_C(x)} \end{bmatrix}$$

Assume that  $C=2$ ,  $y=1$

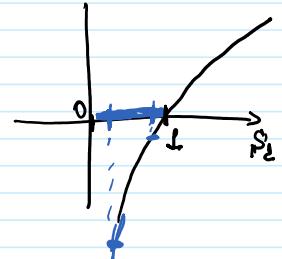
$$\left| \begin{array}{l} \hat{s}_1(x) > \hat{s}_2(x) \Leftrightarrow h_1(x) > h_2(x) \Rightarrow \hat{y}=1 \quad (\text{no error}) \\ \text{Then, } \ell_{\text{CE}}([1, 0], \hat{s}(x)) = -1 \cdot \ln(\hat{s}_1(x)) - 0 \cdot \ln(\hat{s}_2(x)) = -\ln(\hat{s}_1(x)) \approx 0 \end{array} \right.$$

Then,  $\text{ACE}([\underline{0}], \underline{\xi}(x)) = -1 \cdot \ln(f_1(x)) - 0 \cdot \ln(f_2(x)) = -\ln(f_1(x)) \approx 0$

$\Rightarrow f_2(x) > f_1(x) \Leftrightarrow h_2(x) > h_1(x) \Rightarrow \hat{y} = 2$  (mistake)

Then,  $\text{ACE}([\underline{0}], \underline{\xi}(x)) = -1 \cdot \ln(f_1(x)) - 0 \cdot \ln(f_2(x)) = -\ln(f_1(x)) \approx 0$

$$\text{ACE} = \frac{1}{N} \sum_{n=1}^N \text{ACE}(t_n, \xi(x_n)) \leftarrow \text{this is differentiable.}$$



Comment:

If turns out (details omitted for the moment) that minimizing ACE depresses the empirical gl loss (it does not necessarily minimize it)

$$\frac{d\text{ACE}(\tilde{w})}{d\tilde{w}} = \frac{1}{N} \tilde{X}^T (\underline{\xi} - \bar{T}) \in \mathbb{R}^{(D+1) \times G}$$

Gradient matrix.  
 $\in \mathbb{R}^{(D+1) \times G}$

$\begin{bmatrix} \underline{\xi}^T \\ \vdots \\ \underline{\xi}_N^T \end{bmatrix} \in \mathbb{R}^{N \times G}$

$\begin{bmatrix} \bar{T}^T \\ \vdots \\ \bar{T}_N^T \end{bmatrix} \in \mathbb{R}^{N \times G}$

$$\triangleq \begin{bmatrix} \left( \frac{d\text{ACE}}{d\tilde{w}_1} \right)^T \\ \vdots \\ \left( \frac{d\text{ACE}}{d\tilde{w}_G} \right)^T \end{bmatrix}$$

Comment: setting the gradient matrix to 0 and trying to solve for  $\{\tilde{w}_k\}_{k=1}^G$  does not produce a closed-form solution

Only choice left: use a numerical minimization approach:  
Gradient Descent (with backtracking)

$$\begin{aligned}\tilde{W}^{\text{new}} &\leftarrow \tilde{W}^{\text{old}} - \gamma \frac{d \text{ACE}}{d \tilde{W}} \Bigg|_{\tilde{W} = \tilde{W}^{\text{old}}} = \\ &= \tilde{W}^{\text{old}} - \gamma \frac{1}{N} \tilde{X}^T (\tilde{S} - \tilde{T})\end{aligned}$$

↑ Computed using the weights  
in  $\tilde{W}^{\text{old}}$ .

## ② Regarding CE

Divergence measure for distribution

Kullback-Leibler (KL) divergence

$$D_{\text{KL}}(Q \parallel P) \triangleq \int_{x \in \mathbb{R}^D} q(x) \ln \left( \frac{q(x)}{P(x)} \right) dx = \underbrace{\int_{x \in \mathbb{R}^D} q(x) \ln q(x) dx}_{\equiv -\text{Ent}(Q)} - \underbrace{\int_{x \in \mathbb{R}^D} q(x) \ln p(x) dx}_{\equiv \text{CE}(Q, P)}$$

↑  
distributions  
with PDFs  
 $q(x)$  &  $p(x)$

If  $Q \equiv P$  ( $\Leftrightarrow q(x) = p(x)$ ), then  $D_{\text{KL}}(Q \parallel P) = 0$

One can show that  $D_{\text{KL}}(Q \parallel P) \geq 0$

$D_{\text{KL}}(Q \parallel P) \neq D_{\text{KL}}(P \parallel Q)$   
in general

Comments:

$$\text{Ent}(Q) \triangleq - \int_{x \in \mathbb{R}^D} q(x) \ln q(x) dx$$

↓

Convention:  $0 \cdot \ln 0 = 0$

measure of "spread"  
of the distribution

$$\text{CE}(Q, P) = - \int_{x \in \mathbb{R}^D} q(x) \ln p(x) dx$$

$$\text{CE}(Q, Q) = \text{Ent}(Q)$$

$$D_{\text{KL}}(Q \parallel P) = \text{CE}(Q, P) - \text{Ent}(Q) \quad \leftarrow$$

The idea is:

↓ know distribution  
parametrized by  $\theta$

at

The idea is:

$$\underline{\theta}^* = \arg \min_{\underline{\theta}} D_{KL}(Q || P_{\underline{\theta}}) = \arg \min_{\underline{\theta}} CE(Q, P_{\underline{\theta}})$$

unknown distribution  
from which we can  
sample data

know distribution  
parametrized by  $\underline{\theta}$   
 $Q$  &  $P_{\underline{\theta}}$  distribution about some  $\underline{x}$   
or  $(\underline{x}, y)$

$$= \arg \min_{\underline{\theta}} - \int_{\underline{x} \in \mathbb{R}^D} q(\underline{x}) \ln p_{\underline{\theta}}(\underline{x}) d\underline{x} = \arg \min_{\underline{\theta}} \int_{\underline{x} \in \mathbb{R}^D} q(\underline{x}) [-\ln p_{\underline{\theta}}(\underline{x})] d\underline{x}$$

$\rightarrow \mathbb{E}_Q \{-\ln p(\underline{x} | \underline{\theta})\} \approx$

Big takeaway:

Minimizing the neg. log. likelihood  $\equiv CE$

for a chosen model ( $p(\underline{x} | \underline{\theta})$ )  
and data drawn from an unknown  
distribution (which we would like to  
model) is equivalent to minimizing  
the KL divergence of our model  
to the unknown distribution

$$\approx \frac{1}{N} \sum_n [-\ln p(x_n | \underline{\theta})] \propto$$

where  $x_n \sim Q$

$$\propto -\sum_n \ln p(x_n | \underline{\theta})$$

negative log-likelihood  
given the assumed model  
and the data  $\mathcal{D} = \{\underline{x}_n \sim Q\}_{n=1}^N$

The likelihood for MNR classification is

$$\mathcal{L}(\underline{\theta} | \mathcal{D}) = \binom{N}{N_1, \dots, N_G} \prod_{m=1}^N P(y=y_m | \underline{x}_m)$$

typically, ignored

$$\binom{N}{N_1, \dots, N_G} = \frac{N!}{N_1! N_2! \dots N_G!}$$

$\equiv t_{n,k}$

$$\prod_{k=1}^G [P(y=k | \underline{x}_m)]^{I[y_m=k]}$$

$= s_k(x_m)$

$$\Rightarrow -\ln \mathcal{L}(\underline{\theta} | \mathcal{D}) = N \cdot ACE(\underline{\theta})$$

Another big takeaway is: ACE is proportional to the neg.-log-likelihood of the MNR classifier

Relationship between 0/1 loss & CE loss

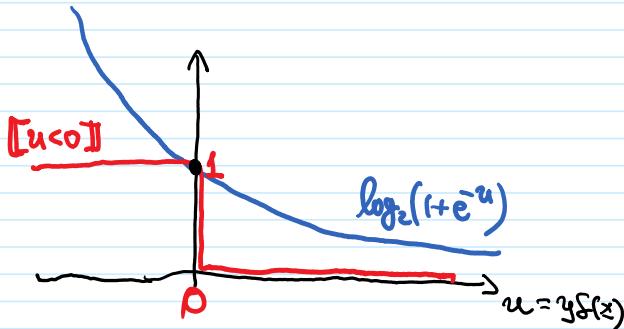
For simplicity, choose  $G=2$

Assume that we have observed  $(x, y)$   $y = \pm 1$  (convenience).

$$\hat{y} = \arg \max_{k=\pm 1} h_k(x)$$

$$\hat{y} = \underbrace{\text{sign}(h_1(x) - h_2(x))}_{S(x) \triangleq} \checkmark$$

$$l_{0/1}(y, \hat{y}) = \llbracket y \neq \hat{y} \rrbracket = \llbracket y\hat{y} < 0 \rrbracket = \llbracket y\delta(x) < 0 \rrbracket \leq \log_2(1 + e^{-y\delta(x)}) =$$



$$= \log_2(1 + e^{-y\delta(x)}) = -\log_2\left(\frac{1}{1 + e^{-y\delta(x)}}\right) \quad \textcircled{1}$$

$$\frac{1}{1 + e^{-y\delta(x)}} = \frac{e^{t_+\delta(x)}}{e^{t_+\delta(x)} + e^{t_-\delta(x)}} = \frac{e^{t_+\delta(x)}}{[1 + e^{\delta(x)}]^{t_++t_-}} = \left(\frac{e^{\delta(x)}}{1 + e^{\delta(x)}}\right)^{t_+} \left(\frac{1}{1 + e^{\delta(x)}}\right)^{t_-}$$

$t_+ = 1 \text{ if } y=1$   
 $t_- = 1 \text{ if } y=-1$

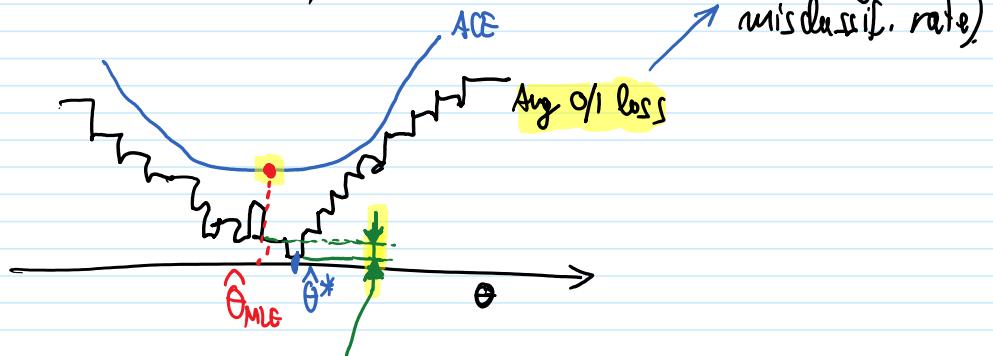
$$= \left(\frac{e^{h_+(x)}}{e^{h_+(x)} + e^{h_-(x)}}\right)^{t_+} \left(\frac{e^{h_-(x)}}{e^{h_+(x)} + e^{h_-(x)}}\right)^{t_-} \quad \textcircled{2}$$

$$= \underline{s}_+(x) \qquad \qquad = \underline{s}_-(x)$$

$$\textcircled{1} \stackrel{\textcircled{2}}{\Rightarrow} l_{0/1}(y, \hat{y}) \leq -t_+ \log_2 \underline{s}_+(x) - t_- \log_2 \underline{s}_-(x) \equiv l_{CE}(t_+, \underline{s}(x))$$

$= t_1 \quad = \underline{s}_1 \quad = t_2 \quad = \underline{s}_2$

Big take away: CE loss upper bounds the 0/1 loss  
(ACE)



hope: that these two  
are "close"

④ Regularization: squared L<sub>2</sub> penalty term.

Useful if #training samples N << dimensionality D of the feature space.

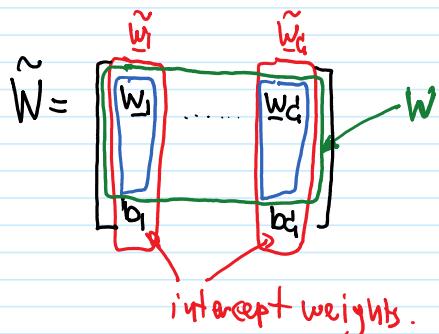
to avoid overfitting.

$$RACE(\tilde{W}) = ACE(\tilde{W}) + \gamma \|W\|_F^2 \quad \text{vs.} \quad \|W\|_F^2 \quad \textcircled{P}$$

penalty parameter  $\gamma \geq 0$

regularized ACE

$$\|W\|_F^2 = \text{trace}\{W^T W\} = \sum_{k=1}^C \|w_k\|_2^2$$



intercept weights  
are excluded

$$\textcircled{P} \Rightarrow \frac{d RACE(\tilde{W})}{d \tilde{W}} = \underbrace{\frac{d ACE(\tilde{W})}{d \tilde{W}}}_{I \tilde{W}^T d - T} + 2\gamma \begin{bmatrix} W \\ 0^T \end{bmatrix} = \frac{1}{N} \left( \tilde{X}^T (S - T) + \gamma \begin{bmatrix} W \\ 0^T \end{bmatrix} \right)$$

$\gamma \rightarrow \frac{1}{2N} \gamma$

$$\bar{N} \perp \tilde{X}^T(S-T)$$

$$\frac{d \|w\|_F^2}{d \tilde{w}} = \underbrace{\begin{bmatrix} \frac{d \|w\|_F^2}{d w} \\ \frac{d \|w\|_F^2}{d b^T} \end{bmatrix}}_{1 \times G} = \begin{bmatrix} 2w \\ 0^T \end{bmatrix}$$

$\downarrow$

$$b \triangleq \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_G \end{bmatrix}$$

Pseudocode:

for various values of  $\gamma \geq 0$  (usually,  $\gamma = 10^e$  where  $e = -4, -3, \dots, 3$ )

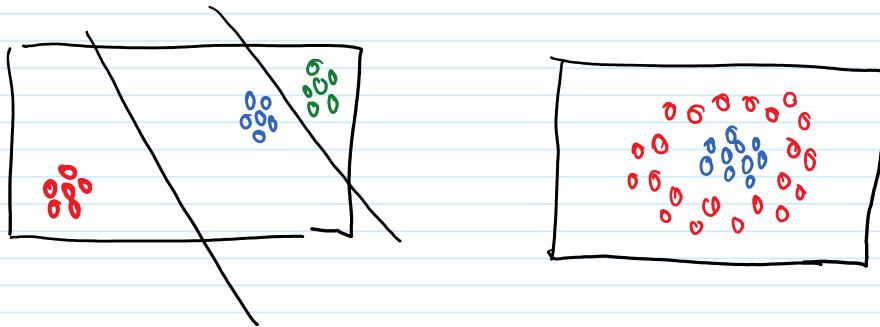
- ▶ Train the MNR classifier for this value of  $\gamma$ . (using training set)
  - ↳ Gradient Descent.

$$\tilde{w}^{\text{new}} \leftarrow \tilde{w}^{\text{old}} - \gamma \left. \frac{d \text{RACE}(\tilde{w})}{d \tilde{w}} \right|_{\tilde{w} = \tilde{w}^{\text{old}}}$$

- ▶ Compute the performance of the resulting MNR classifier on a hold-out (validation) set.

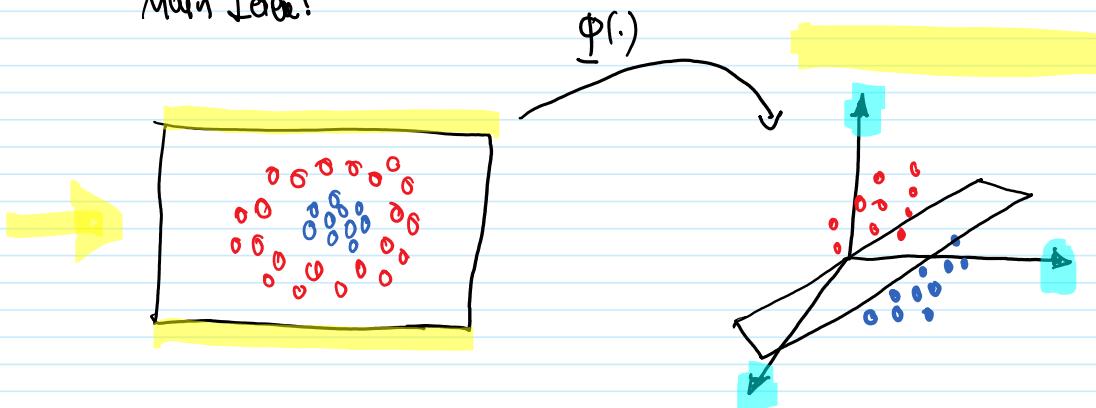
Pick the model ( $\Leftrightarrow$  best value of  $\gamma$ ) that has the best validation performance ("Champion model")

### ④ MNR w/ non-linear transformation of features



For the 2nd example, use a suitable non-linear transformation of the features.

Main Idea:



How to find a suitable  $\underline{\Phi}(\cdot)$ ?

→ Learned from the data

→ "Handcrafted"  $\underline{\Phi}(\cdot)$

→ "Learned"  $\underline{\Phi}(\cdot)$

→ Radial Basis function mapping:

$$\underline{\Phi}(x) = \begin{bmatrix} e^{-\frac{\|x-x_1\|^2}{\xi}} \\ \vdots \\ e^{-\frac{\|x-x_N\|^2}{\xi}} \end{bmatrix} \in \mathbb{R}^N$$

depends on the particular training set.

training samples

where  $\xi$  is a (hyper)parameter

$D_{\text{new}} = N \leftarrow$  may lead to potential overfitting.

Train an MNR classifier using the transformed features  $\{\underline{\Phi}(x_1), \dots, \underline{\Phi}(x_N)\}$ .

