

(1)

# ML lecture 1 - Introduction

$x_n \in \mathbb{R}^{D_x}$ , bold letter  $x$  = tensor

linear model = relationship between  
 $x$  &  $y$

L2 loss = Euclidean = Mean Squared Err

$$\begin{aligned} L_2(w) &= \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2 \\ &= \frac{1}{N} \sum (x_n \cdot w + b - y_n)^2 \\ &= L_2(w) \frac{1}{N} \|xw - y\|^2 \end{aligned}$$

Recall norm  $\|.\|$  = measure size or length of a vector (magnitude)

$\|.\|_{L2}$  = Square, Sum, Square root

Why this sequence of actions?

- Pythag theorem
- Vectors hold direction, summing them plainly may cause  $(1) + (-1)$  to cancel out
- A norm must be non-neg

(2)

Steps for  $\frac{1}{N} \|xw - y\|^2$

① Diff of vectors:  $xw - y$

$$\text{e.g. } [-1-2, 4-1, 3-5] = [-1, 3, -2]$$

vector  $\rightarrow$  vector

② Apply the L2 Norm:  $\|xw - y\|$

- Square Comps  $-1^2, 3^2, -2^2 = [1, 9, 4]$
- Sum  $1 + 9 + 4 = 14$
- Sq Root  $14 = \sqrt{14}$

③ Square result:  $\sqrt{14} \rightarrow 14$

④ Divide by N:  $\frac{1}{N} 14$

Why square at end?

- By definition a norm  $\|\cdot\|$  involves a square root to represent the true distance (Pythag theorem)
- We reapply the square to allow for better/easier maths in future
- To minimize loss we take derivative of the function/equation
- Squared is much easier to work with (chain rule)

(3)

Recall optimal model parameters  
minimize the loss ( $w, b$ )

Training is the process of optimizing  
 $w$  &  $b$

Minimize the loss function:

- find the min when its gradient (derivative) is set to 0
- gradient = Derivative

$$\frac{d L^2(w)}{dw} = 0 \quad \left. \begin{array}{l} \text{deriv of loss} \\ \text{func set to 0} \end{array} \right\}$$

Expand & Sub in loss function:

$$\cancel{\frac{1}{N} \frac{d(\|w \cdot g\|^2)}{dw}}$$

recall  $\frac{d(w)}{dw}$  is the notation that means derivative of a function w/  
respect to w

$D$  = infinite small

$d(w)$  = " change in function

$dw$  = " change in var w

(4)

Keep expanding numerator:

$$\frac{d((xw - y)^T(xw - y))}{dw} = 0$$

$$\hookrightarrow \frac{d(w^T x^T x w - 2w^T x^T y + y^T y)}{dw} = 0$$

Apply derivative w respect to  $w$

$$2x^T x w - 2x^T y$$

Rearrange to get  $w$  (goal)

$$w = (x^T x)^{-1} x^T y$$

$w$  gives optimal linear model

L1 vs L2 Norm:

- ~~L2~~ when differentiated, L2 has no  $w$  term
- L1 gradient is a constant
- constants do show rate of change
- can't see how close a result is

(5)

## Model Generalizability errors: Bias & variance

high bias = less sensitive to train =  
 less complex = risk of under fit

high Sens = More sensitive to train =  
 more complex = risk of overfit

overfitting = memorize training  
 data. Will not generalize well  
 as has not learnt from data

## Regularization:

- To tackle overfitting
- Penalty to loss function
- 

$L_1$  reg = Lasso

$$L_{\text{lasso}}(\omega) = \frac{1}{N} \|\omega - y\|^2 + \alpha \|\omega\|_1$$

reg term

- Penalizes complexity
- Complexity  $\propto$  WS

- L1 Regularization penalize non-zero weights
- encourages  $w=0$  to reduce complexity

## L2 Regularization:

- +  $\alpha \sum w_i^2$
- encourages small weights to reduce complexity

## Classification Models:

$$\{(x_n, y_n)\}_{n=1}^N, x_n \in \mathbb{R}^{D_x}, y_n \in \mathbb{I}^{D_y}$$

Vector  $\rightarrow$  category

$$f(x) = \sigma(xw + b) = \hat{y} \quad \sigma(\cdot) \text{- activation function}$$

Classification sign loss: or 0-1 loss

$$L(w) = \frac{1}{N} \sum I[\text{Sign}(x_n w + b) \neq y_n]$$

$$\text{Hinge loss: } L(w) = \frac{1}{N} \sum \max(0, -y_n(x_n w + b))$$

linear model

- hinge allows how far yes or no

Pos or Neg