

ECE-GY 6143

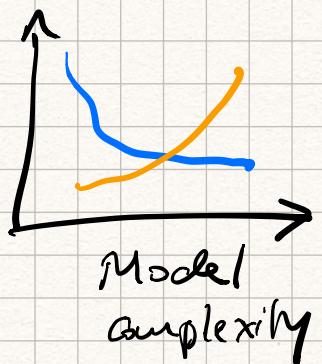
Intro to ML

- ✓ Recap: Model selection
- ✓ Regularization
- ✗ Logistic Regression.

- Bias vs Variance
- Thm's rule

▷ low bias
("higher" complexity)

▷ low variance
("lower" complexity).



Regularization

- Method to control complexity of an ML model.

Components

- Representation
- Loss function
- Optimization

✓
?
✓

Instead of minimizing MSE,
define a new loss function

$$L(\omega) = \text{MSE}(\omega) + \alpha \phi(\omega)$$

scalar ≥ 0
"Regularizer"
Promote simpler ω .
Penalize "unlikely" ω

Examples -

$$1) \quad \phi(\omega) = \|\omega\|_2^2 \quad \begin{matrix} \text{squared L2} \\ \text{-norm.} \end{matrix}$$

$$2) \quad \phi(\omega) = \|\omega\|_1 \quad \text{L1 norm.}$$

$$3) \quad \phi(\omega) = \frac{1}{2} \|\omega\|_2^2 + \frac{\lambda}{2} \|\omega\|_1 \quad \text{"Elastic net".}$$

$$4) \quad \phi(\omega) = \sum_{i=1}^d \left(\omega_i - \frac{1}{d} \sum_{j=1}^d \omega_j \right)^2$$

$$\phi(\omega) = \sum_{i=1}^{d-1} (\omega_i - \bar{\omega})^2$$

Interpretation in terms of linear regression

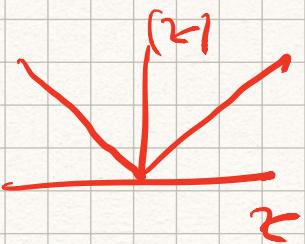
$$① \quad L(\omega) = \frac{1}{2} \|y - X\omega\|_2^2 + \alpha \|\omega\|_2^2$$

Set $\nabla L(w) = 0$ and solve for w^*

$$w^* = (X^T X + \alpha I)^{-1} X^T y.$$

"Ridge regression"

② $L(w) = \frac{1}{2} \|y - Xw\|_2^2 + \alpha \|w\|_1$



$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

cannot optimize using gradient descent
but can use other algorithms

{ sub-gradient descent, LARS,
fixed point continuation, ... }

"LASSO regression"

$$\phi(w) = \|w\|_2^2 \rightarrow \text{Shrinks values of } w^*$$

$$\phi(w) = \|w\|_1 \rightarrow \text{Sparsifies } w^*$$



logistic regression

thus far: regressions

Training Dataset: (x_i, y)
Data \downarrow \uparrow label (continuous)

Q: How to model classification

$$\begin{array}{ll} y = 0 & \rightarrow \text{NO} \\ y = 1 & \rightarrow \text{YES} . \end{array}$$

Pros: Can use linear regression directly.

Cons:

$$y = \langle w, x \rangle$$

- Not scale invariant.
[$100x$ gives $100y$].
- Labelling not unique.

Solution: Instead of finding model f

s.t.

$$y_i = f(x_i)$$

instead we predict

$$\text{Prob}(y_i = 1 | x_i) = f(x_i).$$

Conditional probability of predicting label YES given data point x_i .

$$\Rightarrow \text{Prob}(y_i = 0 | x_i) = 1 - f(x_i). \\ \text{[assuming 2 classes].}$$

Combine into one equation.

$$\text{Prob}(y_i | x_i) = \frac{f(x_i)^{y_i} (1-f(x_i))^{1-y_i}}{\text{Dataset}}$$

Dataset

$$(x_1, y_1)$$

$$(x_2, y_2)$$

:

$$(x_n, y_n)$$

||

$$X = [x_1, \dots, x_n] \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

with independent
samples.

$$\text{Prob}(y | X) = p(y_1 | x_1) p(y_2 | x_2) \dots p(y_n | x_n)$$

$$= \prod_{i=1}^n$$

$$\text{prob}(y_i | x_i)$$

$$= \prod_{i=1}^n f(x_i)^{y_i} (1-f(x_i))^{1-y_i}$$

$\stackrel{i=1}{\text{Likelihood}}$ of training dataset.

More convenient to take \log .

$$-\log \text{Prob}(y|x) = - \left[\sum_{i=1}^n y_i \log f(x_i) + (1-y_i) \log(1-f(x_i)) \right]$$

Negative log likelihood

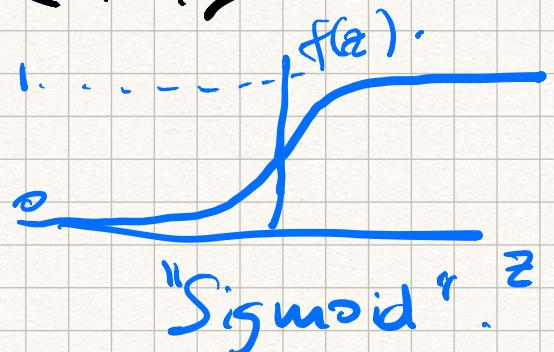
$$\ell_i(f) \left[y_i \log f(x_i) + (1-y_i) \log(1-f(x_i)) \right]$$

"Cross-entropy".

$$L(f) = \sum_{i=1}^n \ell_i(f).$$

$$f(x_i) = \frac{1}{1 + e^{-\langle \omega, x_i \rangle}}$$

$$f(z) = \frac{1}{1 + e^{-z}}$$



$$L(\omega) = \sum_{i=1}^n \ell_i(\omega)$$

$$\ell_i(\omega) = - \left(y_i \log \frac{1}{1 + e^{-\langle \omega, x_i \rangle}} + (1-y_i) \log \frac{1}{1 + e^{\langle \omega, x_i \rangle}} \right)$$

$$1 - \frac{1}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}}$$

$$(1-y_i) \log \frac{e^{-\langle w, x_i \rangle}}{1+e^{-\langle w, x_i \rangle}}.$$

Algorithm to minimize $L(w)$

: gradient descent.

Fact: $\frac{df(z)}{dz} = f(z)(1-f(z))$

GD
 $w_{k+1} = w_k + \alpha_k \sum_{i=1}^n \left(y_i - \frac{1}{1+e^{\langle w, x_i \rangle}} \right) x_i$

Multiclass classification.

k classes $1, \dots, k$.

$$\begin{aligned} \text{Prob}(y_i = k | x_i) &= \frac{\exp(\langle w_k, x_i \rangle)}{Z} \end{aligned}$$

Z : normalization "partition function"

$$Z = \sum_{k=1}^K \exp(\langle w_k, x_i \rangle).$$

$$L(\omega) = \sum_{i=1}^n l_i(\omega)$$

$$l_i(\omega) = \sum_{k=1}^K \mathbb{1}(y_i = k) \log \frac{\exp(\langle \omega_k, x_i \rangle)}{Z}$$

"Softmax"