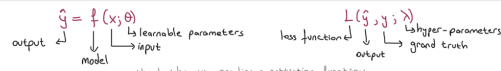


2- ml background



- linear models implement a linear decision boundary (ex: SVM with $w \cdot x$)
- non-parametric models almost everytime implements non-linear decision boundaries (ex: kNN)
- higher capacity \rightarrow more complicated decision boundaries (high complexity can cause overfitting)
- VC dimension = largest number the model can correctly label arbitrary configurations of n points, for binary classification



loss functions:

- $0-1$ = counting how many mistakes $L(w) = \frac{1}{N} \sum_{i=1}^N I(y_i - f(x_i; w) \neq y_i)$ (not used because not differentiable)
- L_2 = $L_2(w) = \frac{1}{2N} \sum_{i=1}^N (y_i - f(x_i; w))^2$ where $y_i = y_i, f(x_i; w) = \sum_{j=1}^J w_j \phi_j(x_i)$ in regression
- hinge = $L_{\text{hinge}}(w) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_i f(x_i; w))$ where $y_i = y_i, f(x_i; w) = \sum_{j=1}^J w_j \phi_j(x_i)$ \rightarrow SVM and max margin alg uses
- log-loss = $L_{\text{log}}(w) = -\frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_i f(x_i; w)})$ where $y_i = y_i, f(x_i; w) = \sum_{j=1}^J w_j \phi_j(x_i)$ similar to hinge but non-zero penalty to all
- cross-entropy = usually used with softmax

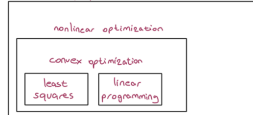
$f(z) = \frac{e^z}{\sum_j e^{z_j}}$ normalized softmax $\rightarrow -\sum_j y_j \log(f(z_j))$

- softmax is not a loss function, it is a normalization function

convexity = a func is convex if the line segment between any two distinct points on the function graph lies above the graph between two points

if the second order derivative is positive, then convex

mathematical optimization



gradient descent = $w^{\text{next}} = w^{\text{curr}} - \eta \nabla_w L(w)$

derivative = rate of change

gradient = vector of partial derivatives $\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$ (used in partial derivative)

subgradient = substitute for gradient when it does not exist

subgradient lines

gradient = the loss is calculated all the examples in the dataset; w updated accordingly

stochastic gradient = set of randomly chosen m example is used to update $w \in [1, N]$

\rightarrow when the learning rate decreases, it converges almost surely

\rightarrow when to stop \rightarrow change in loss or parameters becomes to small

\rightarrow performance in validation set does not improve

\rightarrow max time or iteration is reached

\rightarrow used to prevent overfitting

regularization = controlling the capacity of a learning model \rightarrow lower capacity is favored

\rightarrow margin maximization in SVM \rightarrow choosing best boundary to separate data

$$L(w) = \underbrace{\frac{1}{2} \|w\|^2}_{\text{cost}} + \underbrace{C \sum_{i=1}^N \max(0, 1 - y_i \cdot w \cdot x_i)}_{\text{regularization}} \quad \text{margin} = \frac{2}{\|w\|}$$

\rightarrow dropout and data augmentation is used in regularization

$$L_2 \text{ and } L_1 \text{ norm} = f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \quad \text{simpler} \rightarrow \theta_0 + \theta_1 x_1 \text{ if } \theta_2 \text{ and } \theta_3 \text{ small}$$

$$\text{cost} = \underbrace{\text{training error}}_{\text{loss func}} + \underbrace{\gamma \sum_{i=1}^N \|\theta_i\|}_{\text{L1 regularization}} \rightarrow \text{simpler equations}$$

$$\text{cost} = \underbrace{\text{training error}}_{\text{loss func}} + \underbrace{\lambda \sum_{i=1}^N \|\theta_i\|^2}_{\text{L2 regularization}} \rightarrow \text{simpler equations}$$

model parameters = learned on training set \rightarrow learning / training (with numerical optimizations) (like gradient descent)

hyper-parameters = learned on the validation set \rightarrow model selection / tuning \rightarrow with grid search, or etc

\rightarrow model related = # layers in NN, C in SVM

\rightarrow learning method related = learning rates, batch size..

* cross validation is used on both

$$\text{maximum likelihood (ML)} = \arg \max_{\theta} \prod_{i=1}^N p(x_i | \theta) \rightarrow \arg \max_{\theta} \sum_{i=1}^N \log p(x_i | \theta)$$

to minimize, log both sides

if this is uniform, ML & MAP are equivalent

$$\text{maximum a-posteriori (MAP)} = \arg \max_{\theta} \prod_{i=1}^N p(x_i | \theta) \cdot p(\theta) \rightarrow \arg \max_{\theta} \sum_{i=1}^N (\log p(x_i | \theta) + \log p(\theta))$$

prior

\rightarrow both are used to estimate parameters