

Machine Learning.

(A neural network)

Multilayer Perceptron

$$Q_6(T)$$
 $G(z)$

$$x_{1,1} \cdot w_1 + x_{1,2} \cdot w_2 + \dots + b_0$$

base.

 z_1

$\lambda_{1,1}$ $\lambda_{1,2}$ Input data ①

Complex model.

Multiple filters:

$$w_1, w_2, w_3, \dots$$

W_1, W_2, W_3, \dots, C
Used for non linear seperable
 $Z_i = x_{i,1} \cdot W_{i,1} + \dots$ data.

$$z_i = x_{i,1} \cdot w_{1,1} +$$
$$\underline{x_{1,2} \cdot w_{1,2} + \dots + b_1}$$
$$z_1 = x_1 \odot w_1 + b_1$$
$$\underline{z_2 = x_1 \odot w_2 + b_2}$$

this z value.

$$T = (C \odot Z + t)$$

$$z = w \odot x_1 + b \quad (b_1, b_2, \dots)$$

$$(z_1, z_2, \dots) (w_1, w_2, \dots) \\ = (C \odot (W \odot x_1 + b) + t)$$

$$= F \odot x_1 + f$$

This is same as LR

Scanned with CamScanner

NOTES

There are different models that we can use in ML.

Ex:- Logistic regression, MLP, CNN, ...

Overfitting

Using complex models can perform well on training data and not so well in real (Not generalized model). When we

How to choose a suitable model? increase the
Trying the model in real depth of our model → increase
word by getting new the parameters of the model. →
real word data set increase error rate.

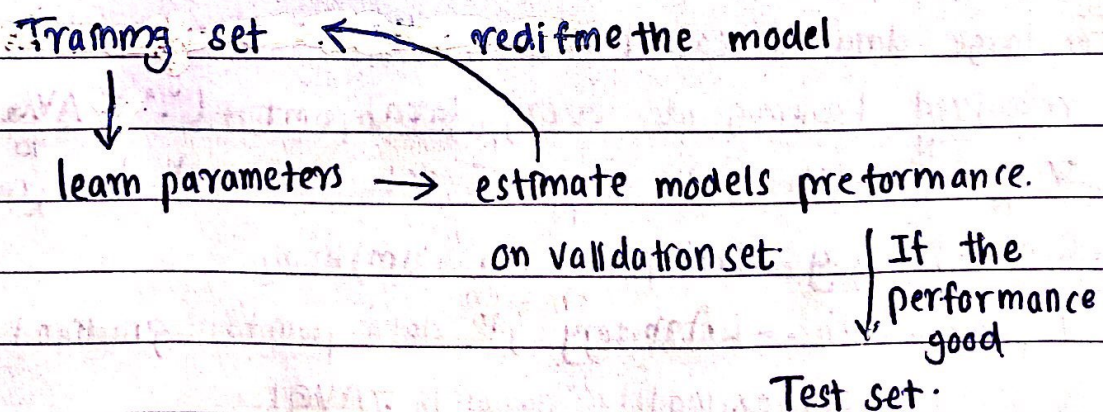
X. Costly. Complex relationships may be too
∴ Split our training complex for reality.
data for 3 groups.

Training → Train the model.

Validation - Compare each approach. Use repeatedly to select best performance
Testing - Evaluate performance of our model by running the
test set.

Equal to running new data (real world data)

Should only be used once. If not lead to bias.



NOTES

How to choose ^(best) optimized parameters?

The values of parameters when the loss function is minimum. — Gradient descent.

$L(x_i, L_i)$ — The distance between the true label and the predicted label by the model for i^{th} input.

We need to take the average loss value.

$$AVG L = \frac{1}{N} \sum_{i=1}^N L(x_i, L_i)$$

All the input.

Normally

$$= -\log P(L_i | \sigma(z_i))$$

$\sigma(z_i)$ — Predicted output value prob

Gradient descent

Take the gradient of

the loss func in an

(update)

arbitrary point. Then move the

θ_{t+1}

Slope

gradient

arbitrary point in

Mostly used cross entropy = $-\sigma(z_i) \log(L_i)$ direction the gradient

We take α

In ML α is the

~~reduces~~ descent (step size)

value when

parameters that

$$\theta_{t+1} = \theta_t - \alpha \nabla f(\theta_t)$$

loss function

we have to learn.

is at it's local minimum.

We keep updating the

parameters until loss func

gradient at

point θ_t

$(\nabla \alpha f)$ is minimum. And the gradient ~~gradient~~ ^{is taken w.r. parameters we have to learn.}

For large data set calculating gradient descent this way

was required looking at every input data point ($\because \nabla AVG L =$

$$\nabla \frac{1}{N} \sum_{i=1}^N L(x_i, L_i) = \frac{1}{N} \sum_{i=1}^N \nabla L(x_i, L_i)$$

gradient of each data point and sum them)

Approximation — Randomly ^{chosen} subset of N or a points (J) gradient

is approximately equal to $\nabla AVG L$

$$\nabla L(x_j, L_j) = \nabla AVG L \text{ or } \nabla \frac{1}{|J|} \sum_{i \in J} L(x_i, L_i)$$

Size of the subset.

NOTES

This is known as Stochastic gradient descent.

Not accurate as gradient descent

But fast. Does more updates within a single time.

May go in both gradient ascending and descending directions. But avg at descending way.

In SGD point update $\theta^{t+1} = \theta^t - \alpha \underbrace{\nabla f(\theta^t)}_{\text{Estimated gradient}} = \nabla f(\hat{\alpha}^t)$

Early stopping.

Optimization goal - Do as well as possible in training set

Validation / Generalization goal - Maximize performance in real world.

Combine optimization with validation process.

Run training data set (with every iteration avg loss ↓)

While running training data we will run validation set also once in a while. Not for parameter θ . To estimate real world performance optimization.

With every iteration compute avg loss. Before we choose optimized parameters when avg loss converge to 0. But now we will stop optimization when validation avg loss stops improving.

Advantage - Save computational cost.

Perform better in real world.

