

BLG 561 E FALL 2021

Deep Learning

19. 10. 2021

Görde ÜNAL

All ML algorithms have 3 components: (model)

① Hypothesis Class: set of functions we consider for mapping input to output

② Loss function: how "good" our selected hypothesis is

③ Optimization: How do we find a hypothesis (\equiv model parameters) w.r. the given loss (objective) fn.

→ Probabilistic view of ML → relates mainly to both hypothesis fn. & the loss fn.

ML Problem. Given $x^{(1)}, \dots, x^{(m)}$: How do I find a probability distribution that captures this data "well"?

$P(\theta; X)$: How do I find θ (parameters) of this distrib. that fits the data well?

→ Maximum Likelihood Estimation (MLE) : the prob. of observing this data :

Given observed independent data points

$$P(x^{(1)}, x^{(2)}, \dots, x^{(m)}; \theta) = \prod_{i=1}^m p(x^{(i)}, \theta)$$

Basic idea in MLE is →

Find the parameters that maximize the probability of the observed data:

$$\max_{\theta} \prod_{i=1}^m p(x^{(i)}, \theta) = \max_{\theta} \underbrace{\ell(\theta)}_{\text{likelihood}} \Rightarrow \text{maximize log likelihood of our data.}$$

$$\rightarrow \max_{\theta} \sum_{i=1}^m \log p(x^{(i)}; \theta) \quad \left(\text{many ML algo. can be interpreted as an MLE} \right)$$

MLE for Linear Regression: $y = \theta^T x + \varepsilon$

Given x , prob. distrib. of y ?

$$\varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$P(y|x; \theta) = \mathcal{N}(\theta^T x, \sigma^2) \propto \exp\left(-\frac{(y - \theta^T x)^2}{2\sigma^2}\right)$$

MLE: $\max \log p(y|x; \theta)$

$$\max_{\theta} \sum_{i=1}^m \log p(y^{(i)}|x^{(i)}; \theta) \rightarrow \min_{\theta} \sum_{i=1}^m \log p(y^{(i)}|x^{(i)}; \theta)$$

(negative log likelihood go to NLL) $\rightarrow \min_{\theta} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$

Least squares loss (MSE) can be viewed MLE under Gaussian assumption

Exercise: Now, say errors are Laplace-distributed; $p(x) = \frac{1}{2\sigma} e^{-|x-\mu|/\sigma}$

Derive MLE in this case \rightarrow (obtain the MAE)

MLE for Logistic Regression: (Binary Classification) $\{+1, -1\}$

$$p(y=+1 | x; \theta) \propto \exp(\theta^T x), \quad : \quad \theta^T x \uparrow.$$

$$p(y=-1 | x; \theta) \propto 1. \quad : \quad \theta^T x \approx 0.$$

$$p(y=+1 | x^i; \theta) = \frac{\exp(\theta^T x^i)}{1 + \exp(\theta^T x^i)}; \quad p(y=-1 | x^i; \theta) = \frac{1}{1 + \exp(\theta^T x^i)}$$

\rightarrow Hypothesis fn $h_\theta(x) = \frac{1}{1 + e^{-y^i \cdot \theta^T x^i}}$: Sigmoid fn. (aka. Logistic fn.)

\rightarrow MLE estimate for θ : $\max_{\theta} \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) = \max_{\theta} \sum_i \log \frac{1}{1 + e^{-y^i \theta^T x^i}}$

$\xrightarrow{\text{convert to min}}$ $\min_{\theta} \sum_i \log(1 + e^{-y^i \theta^T x^{(i)}})$: This is how we obtain the logistic regression loss we've seen last time.

Cross-Entropy Loss: (for Binary Classification)

$$\left. \begin{aligned} p(y=+1|x; \theta) &= h_{\theta}(x) \\ p(y=-1|x; \theta) &= 1 - h_{\theta}(x) \end{aligned} \right\}$$

Recall. Bernoulli r.v.

$$\left. \begin{aligned} p(x=1) &= \phi \\ p(x=-1) &= 1 - \phi \end{aligned} \right\} p(x) = \phi^x (1 - \phi)^{1-x}$$

Combine $p(y|x; \theta) = (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$.

$$\begin{aligned} \rightarrow \ell(\theta) &= p(y|x; \theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^m (h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}) \end{aligned}$$

we assume indep. data samples

→ Maximize the log-likelihood:

$$\max \rightarrow \ell(\theta) = \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))$$

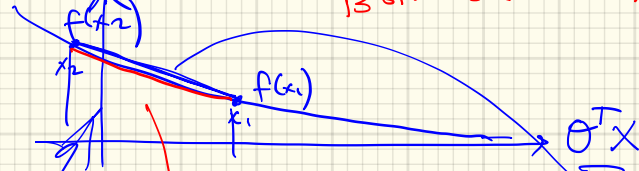
$H(p) = -\sum p \log p$

aka Binary Cross Entropy (BCE)

$$\min \mathcal{L}_{BCE} = \sum_{i=1}^m \underbrace{-y^{(i)} \log\left(\frac{1}{1 + e^{-\theta^T x^{(i)}}}\right)}_{y=+1} - \underbrace{(1 - y^{(i)}) \log\left(\frac{1}{1 + e^{-\theta^T x^{(i)}}}\right)}_{y=-1}$$

Q: Is BCE a convex fn?

for $y=+1$



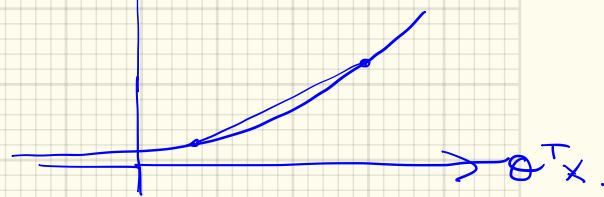
is this convex?

Convex fn:
 $\alpha \in (0,1)$

$$f(\alpha x_1 + (1-\alpha)x_2)$$

Both are convex.

$y=-1$



is this convex?

$$f(\alpha x_1 + (1-\alpha)x_2) \leq \alpha f(x_1) + (1-\alpha)f(x_2)$$

BCE: is sum of 2 convex fn.s. convex? ✓ Yes.

SOFTMAX Loss Fn. Now y takes values $\{1, \dots, k\}$ $y \in \{1, \dots, k\}$

$p(y=j | x; \theta) \propto \exp(\theta_j^T x)$ (unnormalized prob) ← likelihood fn.

Normalize

$$p(y=j | x; \theta) = \frac{\exp(\theta_j^T x)}{\sum_{l=1}^k \exp(\theta_l^T x)}$$

- θ_1
- θ_2
- \vdots
- θ_k

→ maximize log-likelihood.

$$\max_{\theta} \log p(y=j|x; \theta) = \theta_j^T x - \log \sum_{l=1}^k \exp(\theta_l^T x)$$

Softmax loss \Rightarrow $\min_{\theta} \sum_{i=1}^m \log \sum_{l=1}^k \exp(\theta_l^T x^{(i)}) - \theta_{y^{(i)}}^T x^{(i)}$ true class label: $y^{(i)}$.

Multi-class classification softmax loss

$$\min_{\theta} \left[\sum_{i=1}^m \log \sum_{l=1}^k \exp(h_{\theta_l}(x^{(i)})) - h_{\theta_{y^{(i)}}}(x^{(i)}) \right]$$

Multi-Class SVM Loss: $l_{\text{hinge}}^{(i)} = \sum_{j \neq y_i} \max(0, 1 + s_j - s_{y_i})$

Correct class score.

s_j : prediction hypothesis.

if $s_j + 1 \leq s_{y^{(i)}}$ → No loss is added.

When $s_j + 1 \geq s_{y_i}$

+ margin of 1

$S = \theta^T x$: linear hypothesis

the loss (penalty) is added to the loss.

→ minimize overall SVM loss

$$\theta^* = \arg \min_{\theta} L_{\text{svm}}(\theta) = \frac{1}{m} \sum_{i=1}^m l_{\text{hinge}}^{(i)}$$

Q. $L_{\text{svm}}(\theta^*) \rightarrow$ is θ^* unique?

Consider $\sum \max(0, \theta_j^T x^i - \theta_{y^i}^T x^i + 1) \rightarrow \theta$ was a minimizer

Double θ $\left\{ \begin{array}{l} \sum \max(0, 2\theta_j^T x^i - 2\theta_{y^i}^T x^i + 1) \rightarrow 2\theta$ is also a minimizer.

$2\theta_j^T x^i + 1 \geq 2\theta_{y^i}^T x^i$ ✓ How to prevent that!

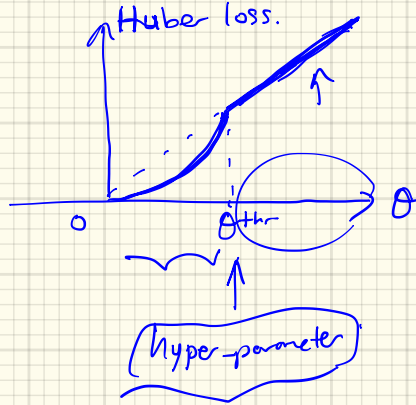
Generally:
 $\arg \min_{\theta}$

$$L(\theta) = \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m l^{(i)}(\text{ho}(x^i), y^{(i)}) + \lambda \cdot R(\theta)$$

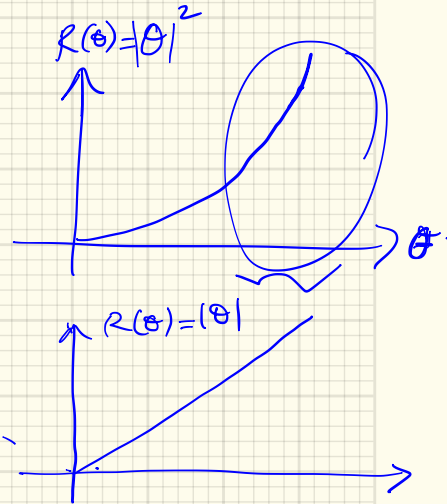
Ex: Regularizer $R(\theta) = \|\underline{\theta}\|_2^2 = \sum_i \theta_i^2$: penalize growth of θ
 L2-constraint Data Term. Add a Regularizer

L1 constraint: $R(\theta) = \|\theta\|_1 = \sum_i |\theta_i|$

L2 + L1
(regularizer)
Huber Loss



Lasso Elastic X.



→ Last time: predicting electricity demand; we used linear hypothesis

Linear Hypothesis Class w/ Nonlinear Features.

Peak electricity demand⁽ⁱ⁾ vs. \downarrow mean Temp.⁽ⁱ⁾ $\{x^{(i)}, y^{(i)}\}$
of a given day.

$$x^{(i)} = \begin{bmatrix} (\text{mean Temp})^2^{(i)} \\ (\text{mean Temp}) \\ 1 \end{bmatrix} = \begin{bmatrix} x_2^{(i)} \\ x_1^{(i)} \\ 1 \end{bmatrix} = \underline{x} = \begin{bmatrix} x_2 \\ x_1 \\ 1 \end{bmatrix}$$

Same hypothesis class as before:

$$h_{\theta}(x) = \underline{\theta}^T \underline{x} = \theta_1 x_2^{(i)} + \theta_2 x_1^{(i)} + \theta_3 : \text{eg. we used a quadratic polynomial}$$

Generalize: n^{th} order polynomial on the data

$$h_{\theta}(x) = \sum_{j=0}^{10} \theta_j x_j^{(i)} : \text{eg. a 10th degree polynomial.}$$

→ Next: check the Slides →

Until now, we've seen Linear Hypothesis fns:

$$\rightarrow y = Wx + b$$

→ Nonlinear Hypothesis Examples: in ML:

$$y = \theta_{n \times 1}^T x_{n \times 1}$$

— Nearest-Neighbor Classifier: predict output based on "nearest" ex. in the training set. $\{x_i\}_{i=1}^m$

$$h(x) = y \left\{ \arg \min_{i \in \{1, \dots, m\}} \text{"dist"}(x, x_i) \right\}$$

Non-parametric hypothesis class. pros: no training required.

con: keep X go over all the training data even at TESTING time:

exercise: Write the hypothesis class for KNN (K-nearest neighbor classifier).



refer to a specific class of

→ Neural Networks: (Nonlinear) Hypothesis Functions.

1 refer to the 1st element. (1)

2 Any loss fn. ✓

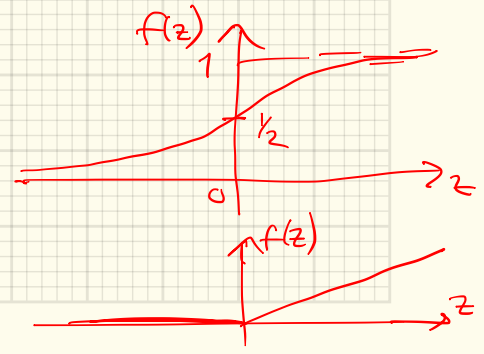
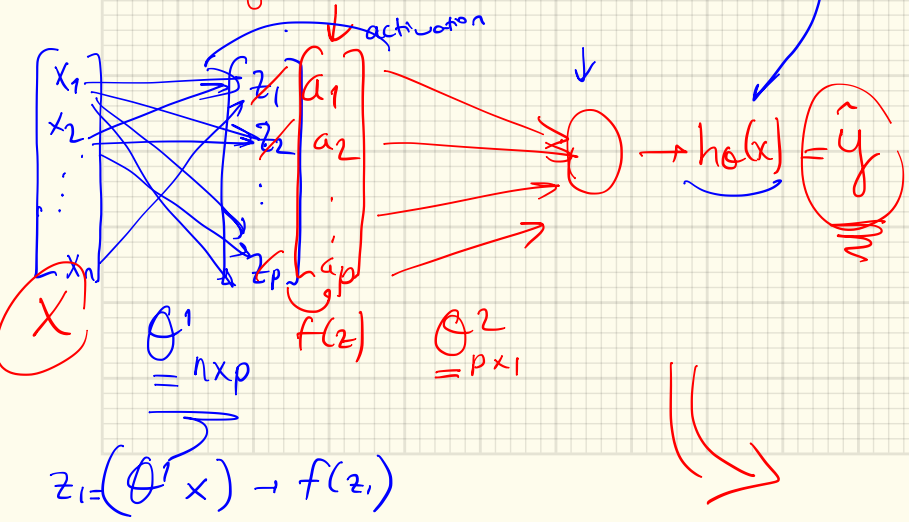
3 Any optimization method ✓

NN hypothesis

Let's consider a NN: $h_{\theta}(x) = \sigma(\theta_2^T f(\theta_1^T x))$

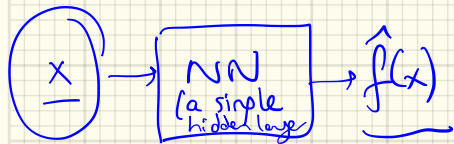
the hypothesis class for a 2-layer NN
 ≡ a single hidden layer NN

↑ ↑ nonlinear
 elementwise $f: \mathbb{R} \rightarrow \mathbb{R}$:
 nonlinearity $f(z) = \frac{1}{1 + e^{-z}}$



Universal Function Approximation Theorem:

Given any smooth fn. $f: \mathbb{R} \rightarrow \mathbb{R}$ (closed region $X \subset \mathbb{R}^n$)
 $\mathbb{R}^n \rightarrow \mathbb{R}$
and $\epsilon > 0$, we can construct a one hidden-layer neural network, \hat{f} s.t.



$$\max_{x \in X} |f(x) - \hat{f}(x)| \leq \epsilon$$

★ A NN w/ a single hidden layer (K enough hidden units)
w/ a squashing (activation) fn. is a universal fn. approximator.

Food for thought: →

It (though a 2 layer NN can represent any function, it may fail due to several reasons: What are they?

→ Importance of the nonlinear activation fn.

Consider a 2 stage hypothesis class.

$$h_{\theta}(x) = \underline{w}_2 \left(\underbrace{\underline{w}_1 x + b_1}_{a_1 \in \mathbb{R}^{k \times 1}} \right) + b_2, \quad x \in \mathbb{R}^n$$

$$\underline{\Theta} = \left\{ \underline{w}_1 \in \mathbb{R}^{k \times n}, b_1 \in \mathbb{R}^k, \underline{w}_2 \in \mathbb{R}^{1 \times k}, b_2 \in \mathbb{R} \right\}$$

Final hypothesis :

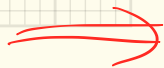
$$h_{\theta}(x) = \underline{w}_2 \underline{w}_1 x + \underline{w}_2 b_1 + b_2$$

$$h_{\theta}(x) = \underline{\tilde{w}} x + \underline{\tilde{b}}$$

The composed model is still affine (linear)

∴ We have not gained any representation power.

$h_{\theta}(x)$: # parameters in $\underline{\Theta} \equiv$ capacity of the model.



NNs introduce a nonlinear fi. after each affine operation.

$$h_0(x) = f_2(\underline{w}_2 f_1(w_1 x + b_1) + b_2)$$

$f_1, f_2 : \mathbb{R} \rightarrow \mathbb{R}$ nonlinear (elementwise) functions

$$f(\underline{x}) = \begin{bmatrix} f(x_1) \\ f(x_2) \\ \vdots \\ f(x_n) \end{bmatrix}$$

Ex: Learning the XOR Function:

Training data $x^{(i)} = \begin{cases} 0, 0 \\ 0, 1 \\ 1, 0 \\ 1, 1 \end{cases}$, $y^{(i)} = \begin{cases} 0 \\ 1 \\ 1 \\ 0 \end{cases}$
 $m=4$

Try a linear hypothesis: $h_0(x^i) = \underline{\theta}^T x + b$

Q: Could a linear work? $y^i = h_0(x^i) = \underbrace{[\theta_1 \ \theta_2 \ b]}_{\underline{\theta}} \begin{bmatrix} x_1^i \\ x_2^i \\ 1 \end{bmatrix}$

② Let's use MSE: $L(\underline{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_0(x^i) - y^i)^2 = (\underline{\theta}^T \underline{x} - \underline{y})^2$

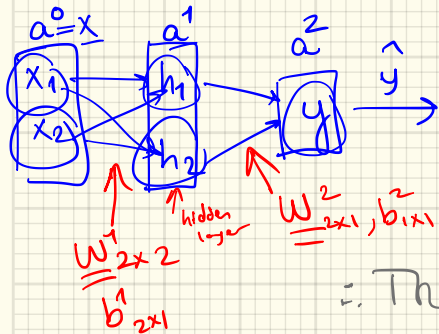
→ Normal eqs. $\underline{\theta} = (\underline{X}^T \underline{X})^{-1} (\underline{X}^T \underline{y}) \rightarrow \underline{\theta}^T = [0 \ 0 \ 0.5]$

→ Linear classifier.

$$h_0(x) = \begin{bmatrix} 0.5 & 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix} = 0.5$$

Linear classifier
not able to learn
represent XOR fn.

→ Use a NN w/ a single hidden layer



$$h = a^1_{2 \times 1} = f(\underline{W}^1 a^0 + \underline{b}^1)$$

nonlinearity affine layer op.

$$\hat{y} = a^2 = \underline{W}^2 a^1 + \underline{b}^2$$

affine layer f(x) sp. f(x) = EU

∴ The Nonlinear hypothesis fn.
for the overall network:



$$\hat{y} = \underline{W}^2 f(\underline{W}^1 a^0 + \underline{b}^1) + \underline{b}^2 = h_0(x), \underline{\theta} = (\underline{W}^1, \underline{b}^1, \underline{W}^2, \underline{b}^2)$$

Suppose we know specify

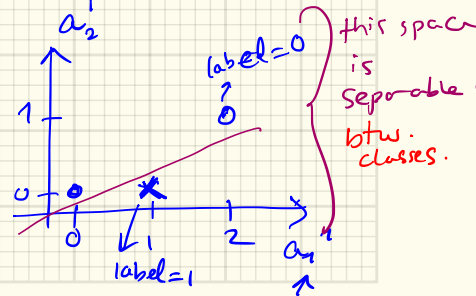
the role to be:

$$\underline{W}^1 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \underline{b}^1 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$\underline{W}^2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \underline{b}^2 = 0$$

Let's check the hidden layer:
for 4 training samples

$$\underline{a}^1 = \begin{cases} a^1(1) \\ a^1(2) \\ a^1(3) \\ a^1(4) \end{cases} = \begin{pmatrix} 0, 0 \\ 1, 0 \\ 1, 0 \\ 2, 1 \end{pmatrix} ; \hat{y} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$



→ ✓ Now, while we increased the capacity of the model \equiv (hypothesis) we increased our representation power.

I learned (the feature space) to represent Xor function.

Also, we fit to the training data well.

Note: Generalization is another story.