


Linear Least Squares

$$J(\theta) = \frac{1}{n} \sum_{i=0}^n (\hat{y}_i - y_i)^2 = \frac{1}{n} \sum_{i=0}^n (x_i \theta - y_i)^2$$

$$\hat{y}_i = x_i \theta$$

matrix notation: $J(\theta) = (X\theta - y)^T (X\theta - y)$

$$\theta = (X^T X)^{-1} X^T y$$

Sigmoid

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

Logistic Regression

$$\hat{y}_i = f(x_i, \theta)$$

$$p(y | X, \theta) = \hat{y} = p(y=1 | X, \theta) \cdot \prod_{i=1}^n p(y_i=1 | x_i, \theta)$$

$$\hat{y} = \prod_{i=1}^n \hat{y}_i \cdot (1 - \hat{y}_i)^{1-y_i}$$

\rightarrow BCE (binary cross entropy)

$$\boxed{\sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)}$$

$$\mathcal{L}(\hat{y}_i, y_i) = y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$

$$\text{of cases} \quad \begin{cases} y_i = 1 & \rightarrow \mathcal{L}(\hat{y}_i, 1) = \log \hat{y}_i \\ y_i = 0 & \rightarrow \mathcal{L}(\hat{y}_i, 0) = \log (1-\hat{y}_i) \end{cases}$$

$$\text{minimization} \rightarrow -\frac{1}{n} \sum_{i=0}^n \mathcal{L}(\hat{y}_i, y_i)$$

$$\text{BCE general case} \rightarrow \mathcal{L}(\hat{y}_i, y_i) = \sum_{j=1}^N \sum_{i=1}^C y_{ij} \log \hat{y}_{ij}$$

Activation Functions

- Sigmoid $\rightarrow \sigma(x) = \frac{1}{1+e^{-x}}$

- tanh $\rightarrow \tanh(x)$

- ReLU $\rightarrow \max(0, x)$

- Leaky ReLU $\rightarrow \max\{0.1x, x\}$

- Parametric ReLU $\rightarrow \max(a x, x)$

L1 Loss

$$L(y, \hat{y}, \theta) = \frac{1}{n} \sum_{i=1}^n \|y_i - \hat{y}_i\|_1$$

MSE (Mean-Squared-Error)

$$\|y - \hat{y}\|_2^2 = (y_i - \hat{y}_i)^2$$

$$L(y, \hat{y}, \theta) = \frac{1}{n} \sum_{i=1}^n \|y_i - \hat{y}_i\|_2^2$$

BCE (Binary Cross Entropy) \rightarrow binary classification

$$L(y, \hat{y}, \theta) = -\frac{1}{n} \sum_{i=1}^n y_i \cdot \log \hat{y}_i + (1-y_i) \cdot \log (1-\hat{y}_i)$$

Cross Entropy \rightarrow multi-class classification

$$L(y, \hat{y}, \theta) = - \sum_{i=1}^n \sum_{k=1}^K (y_{ik} \cdot \log \hat{y}_{ik})$$

Stochastic Gradient Descent (SGD)

$$\frac{1}{n} \left(\sum_{i=1}^n L_i(\theta, x_i, y_i) \right) = E_{i \sim [1 \dots n]} [L_i(\theta, x_i, y_i)]$$

Bias \rightarrow Difference between average prediction of our model and the correct value which we are trying to predict

High error on train and test data UNDERFITTING

High Variance \rightarrow OVERFITTING

Multiclass classification: Softmax

$$p(y_i | x_i, \theta) = \frac{e^{s_{y_i}}}{\sum_{k=1}^c e^{s_k}} = \frac{e^{x_i^\top \theta_{y_i}}}{\sum_{k=1}^c e^{x_i^\top \theta_k}}$$

- Cross Entropy Loss: $L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_k e^{s_k}} \right) \rightarrow$ Always want to improve (loss never going to)

- Hinge loss: $L_i = \sum_{k \neq y_i} \max(0, s_k - s_{y_i} + 1) \rightarrow$ saturates

Xavier Initialization

- Gaussian with 0 mean, with standard deviation:

$$\text{Var}(s) = \sum_i^n \text{Var}(x_i) \text{Var}(w_i) = n [\text{Var}(w) \text{Var}(x)]$$

Ensure the variance of the output is the same as the input

$$\text{var}(s) = \text{var}(x) \rightarrow n \cdot \underbrace{\text{Var}(w) \text{Var}(x)}_1 = \text{Var}(x)$$

$$\rightarrow \text{Var}(w) = \frac{1}{n}$$

- With ReLU (that kills half of the data)

$$\text{Var}(w) = \frac{1}{n/2} = \frac{2}{n}$$

CNN

- Without padding : (Valid convolution)

Input : $N \times N$

Filter : $F \times F$

Stride : 1

$$\text{Output} = \left(\frac{N-F}{S} + 1 \right) \times \left(\frac{N-F}{S} + 1 \right)$$

- With padding : (Same convolution \rightarrow output = input size)

Padding : P

↓
setting
 $P = \frac{F-1}{2}$

$$\text{Output} = \left(\left\lfloor \frac{N+2P-F}{S} \right\rfloor + 1 \right) \times \left(\left\lfloor \frac{N+2P-F}{S} \right\rfloor + 1 \right)$$

Pooling

- Spatial filter $F \times F$

- Stride S

- Output $\rightarrow W_{out} \times H_{out} \times D_{out}$

$$W_{out} = \frac{W_{in}-F}{S} + 1$$

$$H_{out} = \frac{H_{in}-F}{S} + 1$$

$$D_{out} = D_{in}$$

CNN Architectures

- LeNet → Digit recognition (10 classes)

Conv → avg pool → Conv → avg pool → FC → FC → \hat{y}

tanh / Sigmoid as activation function

60K parameters

- AlexNet (Same Conv) (1000 classes)

Conv → max pool → Same conv. → Max pool

→ Same conv → Same conv → Same conv → Max pool → FC = FC → FC

ReLU

~1000 times bigger than LeNet (softmax)

60M parameters

- VGGNet

- Striving for simplicity

- Conv: 3x3 filters with stride 1 (Same conv) } just repeating
Max pool: 2x2 with stride 2 }

. FC , FC , Softmax

137 M parameters

- ResNet

- Xavier / 2 Init.
- SGD + Momentum
- Learning rate 0.1 / 10 when plateau
- Mini-batch 256
- Weight decay $1e^{-5}$
- No dropout

60M parameters

deeper networks ResNet → better performance

- Xception Net (extreme version of Inception)

- 36 conv layers (with skip connections)

- Applies Depthwise Separable Convolutions instead of normal conv

- U-Net

• Left Side → Contraction Path [Encoder]

• Right Side → Expansion Path [Decoder]

LSTM

• Gates:

$$1. \text{ Forget gate } f_t = \text{sigmoid}(\Theta_{xf} X_t + \Theta_{hf} h_{t-1} + b_f)$$

$0 \rightarrow$ forget info (erase cell)
 $1 \rightarrow$ keep info (keep cell)

$$2. \text{ Input gate } i_t = \text{sigmoid}(\Theta_{xi} X_t + \Theta_{hi} h_{t-1} + b_i)$$

decides which values will be updated

$$\tanh(-1, 1)$$

add subtract

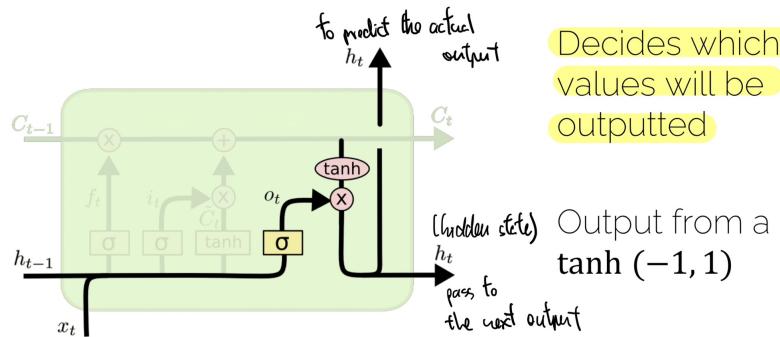
$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

3. Output gate \rightarrow Decides which values will

$$h_t = o_t \odot \tanh(C_t)$$

be outputted

- Output gate $h_t = o_t \odot \tanh(C_t)$



- Forget gate
- Input gate
- Output gate
- Cell update
- Cell
- Output

$$\begin{aligned}
 f_t &= \text{sigm}(\theta_{xf}x_t + \theta_{hf}\mathbf{h}_{t-1} + \mathbf{b}_f) \\
 i_t &= \text{sigm}(\theta_{xi}x_t + \theta_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i) \\
 o_t &= \text{sigm}(\theta_{xo}x_t + \theta_{ho}\mathbf{h}_{t-1} + \mathbf{b}_o) \\
 g_t &= \tanh(\theta_{xg}x_t + \theta_{hg}\mathbf{h}_{t-1} + \mathbf{b}_g)
 \end{aligned}$$

exactly the same in each time step

$$\mathbf{C}_t = \underbrace{f_t \odot \mathbf{C}_{t-1}}_{\text{weights}} + i_t \odot \mathbf{g}_t$$

$$\mathbf{h}_t = o_t \odot \tanh(\mathbf{C}_t)$$

Learned through backpropagation

I2DL · Draft Date

6.1

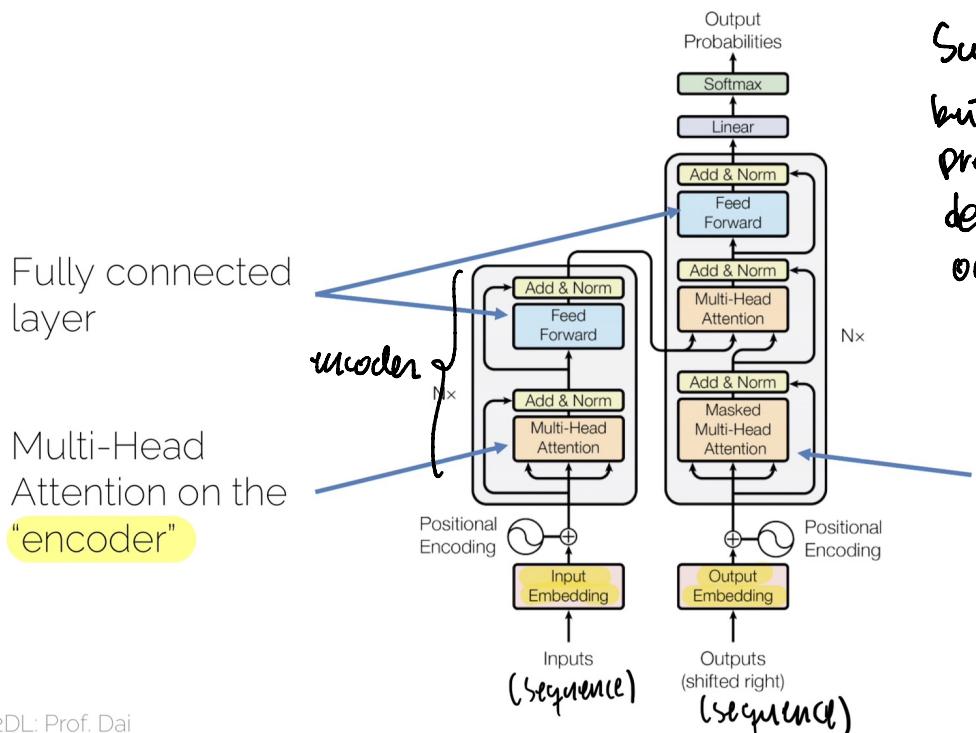
\nearrow f_t for important info \rightarrow so not vanishing gradient

$$\mathbf{C}_t = f_t \odot \mathbf{C}_{t-1} + i_t \odot \mathbf{g}_t$$

\searrow weights

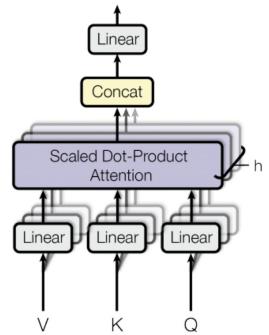
Transformers

Transformers



I2DL: Prof. Dai

Multi-Head Attention



Intuition: Take the query Q , find the most similar key K , and then find the value V that corresponds to the key.

In other words, learn V, K, Q where:

- V – here is a bunch of interesting things.
- K – here is how we can index some things.
- Q – I would like to know this interesting thing.

Loosely connected to Neural Turing Machines (Graves et al.).

Multiple queries with Keys (similarities) get the values

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

"normalize"

Self-attention: complexity

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

where n is the sequence length, d is the representation dimension.

k is the convolutional kernel size, and r is the size of the neighborhood.

Considering that most sentences have a smaller dimension than the representation dimension (in the paper, it is 512), self-attention is very efficient.

What is missing from self-attention?

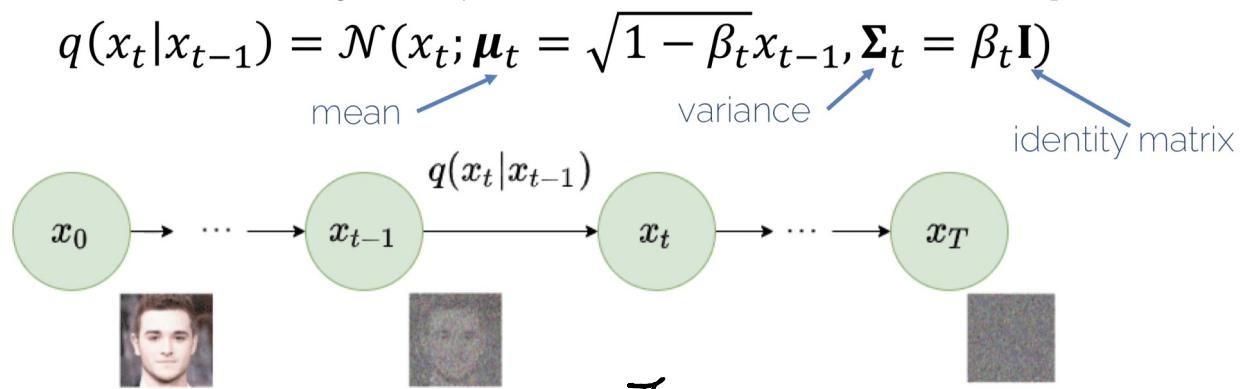
- Convolution: a different linear transformation for each relative position. Allows you to distinguish what information came from where.
- Self-attention: a weighted average.



Diffusion:

1. Forward Diffusion

- Markov chain of T steps
 - Each step depends only on previous
- Adds noise to x_0 sampled from real distribution $q(x)$



$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$$

Not efficient

Reparameterization

- Define $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=0}^t \alpha_s$, $\epsilon_0, \dots, \epsilon_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\begin{aligned} x_t &= \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_{t-1} \\ &= \sqrt{\alpha_t}x_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-2} \\ &= \dots \\ &= \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_0 \end{aligned}$$

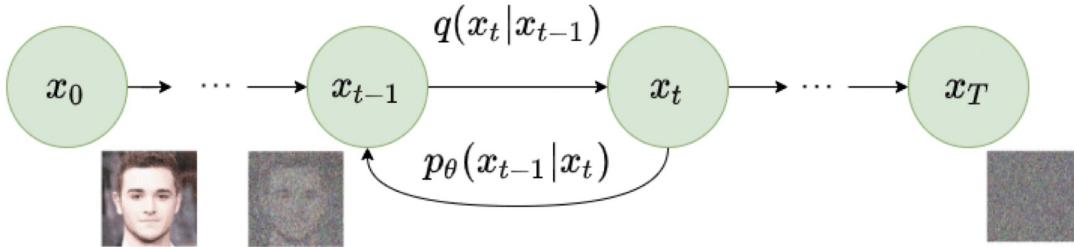
$$x_t \sim q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

2. Approximate Reverse Process

- Approximate $q(x_{t-1}|x_t)$ with parameterized model p_θ (e.g., deep network)

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

$$p_\theta(x_{0:T}) = p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)$$



Loss: Training a Diffusion Model

- Optimize negative log-likelihood of training data

$$\begin{aligned} L_{VLB} &= \mathbb{E}_q \left[D_{KL}(q(x_T|x_0) || p_\theta(x_T)) \right]_{L_T} \\ &+ \sum_{t=2}^T \underbrace{D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))}_{L_{t-1}} - \log p_\theta(x_0|x_1) \end{aligned}$$

- Input and output dimensions must match
- Highly flexible to architecture design
- Commonly implemented with U-net architecture