

# Neural Tangent Kernels (NTK)

★ Looking back at classic kernels

$$x_i \in \mathbb{R}^d \rightarrow \phi(x_i) \in \mathbb{R}^D \quad D \gg d$$

Non-linear

Transforming our features to a higher dimension ( $D$ ) using  $\phi$

Ex-  $d=3$   $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_1 x_2 \\ x_2 x_3 \\ \vdots \end{bmatrix}$

Now,  $f(w, x) = w^T \phi(x)$

weight

features

Transformed  
features

Linear combination

( ignore bias for simplicity )

This model is linear in  $w$  but not in  $x$  (due to non-linear  $\phi$ )

{ we will still have a convex }  
optimisation problem

Gradient descent

$$w(t+1) = w(t) - \eta \nabla_{\underline{w(t)}} L(w(t))$$

$$\sum_{i=1}^n (y_i - f(w, x_i)) \nabla_w f(w(t), x_i)$$

(x<sub>i</sub>)

Static(fixed)

→ As  $d(w) = \frac{1}{2} (y_i - \underbrace{f(w, x_i)}_{\hat{y}_i})^2$

Limitations

1.  $\phi$  is fixed

2.  $\phi(x) \in \mathbb{R}^D$  ( $D \gg d$ )

{ High computational cost }

☆ kernel trick

→ Usually we don't have to worry about computing  $\phi(x)$ s, rather we look at —

$$\langle \phi(x_i), \phi(x_j) \rangle$$

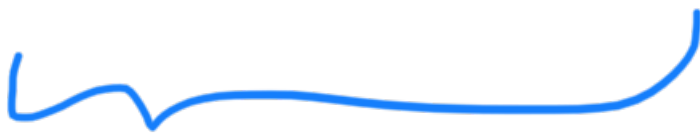
$\langle \rangle \rightarrow$  inner product  
(generalisation of  
dot product)

$$V = \begin{bmatrix} \quad \end{bmatrix}, U = \begin{bmatrix} \quad \end{bmatrix}$$

$$\langle V, U \rangle = V^T U \rightarrow \text{Scalar quantity}$$



$$K(x_i, x_j) \triangleq \langle \phi(x_i), \phi(x_j) \rangle$$



kernel value b/w  $x_i$  &  $x_j$

$$\underbrace{K}_{\text{Kernel matrix}} \in \mathbb{R}^{m \times n} = \begin{pmatrix} & & & \\ & & & \\ & & & \\ & & & \end{pmatrix} \geq 0$$

Kernel matrix

$(i, j)$  element  
corresponds to  
 $k(x_i, x_j)$



$$K \geq 0$$

&  $K$  is symmetric

usually can be computed without  
"explicitly" computing  $\Phi(x)s$

Example -

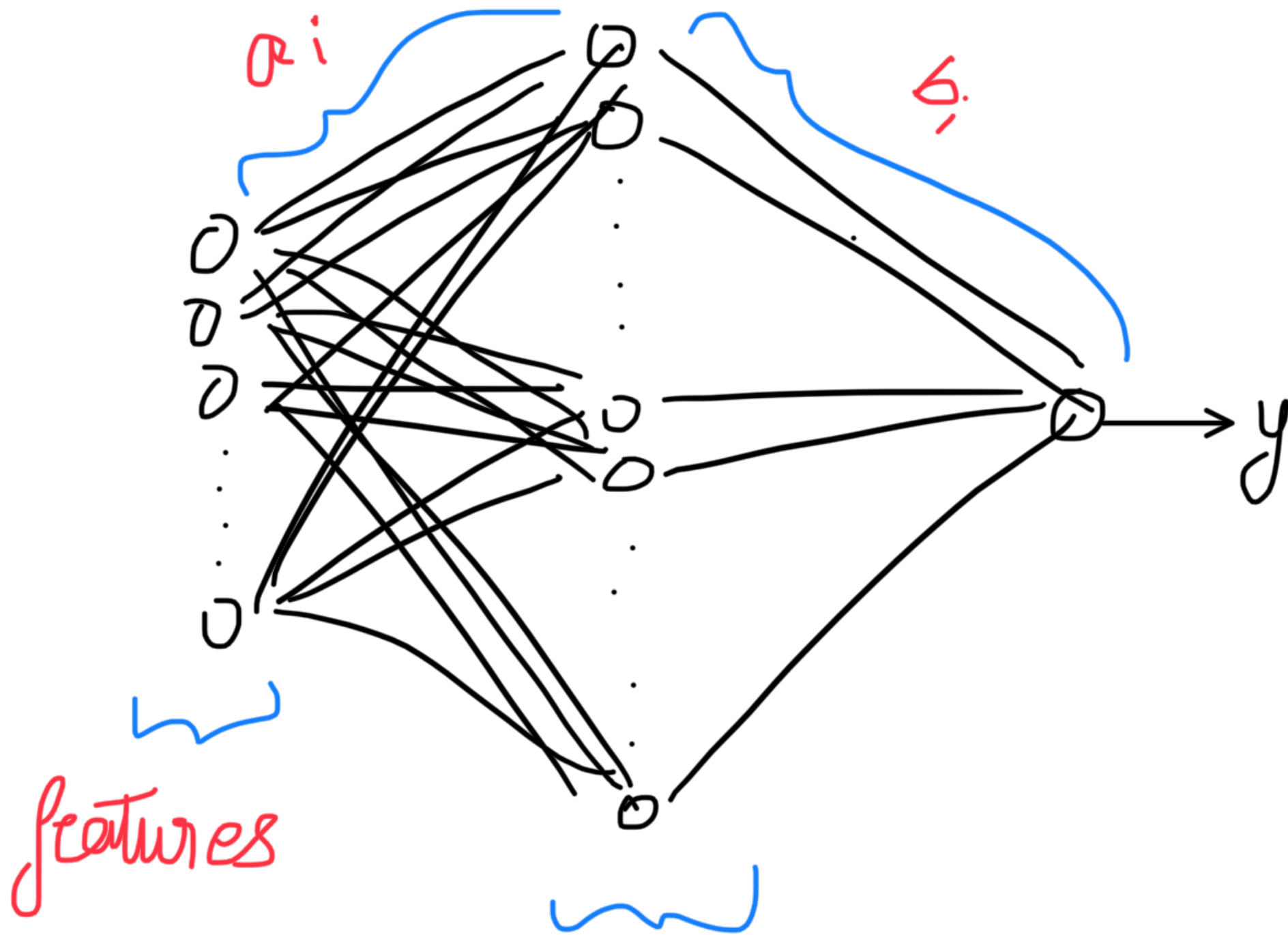
1. Polynomial Kernel

$$K(x_i, x_j) = (c + x_i^T x_j)^K$$

Nonlinear Transform Kernel  $(x_1 + 1, \dots)$



Neural Network Models (NN / KS)



only 1 hidden layer with

$m$  neurons

$$\text{let } nn = y = f(w, x)$$

for simplicity ignore bias terms

$$y = f(w, x) = \sum_{i=1}^m b_i \sigma(a_i^T x)$$

$$y = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \sigma(a_i^T x)$$

Ahead (why?)

Now, taking SSE Loss

$$L(w) = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

Minimising  $L(w)$  using Gradient descent

$$w(t+1) = w(t) - \eta \nabla L$$

Dynamic ( $w(t)$ )

{ Static in Linear  
regression }

$$\nabla \mathcal{L} = \sum_{i=1}^n (f(w, x_i) - y_i) \nabla_w f(w(t), x_i)$$

depends on  $w(t)$   
hence dynamic

lets initialize weights using

Normal distribution —

$$W_{init} \sim N(0, 1)$$

$$W_0 \rightarrow W_1 \rightarrow W_2 \rightarrow \dots W_n$$

Gradient descent

when  $m \rightarrow$  very large

width of neural net

— unlike neural network

$w_0, w_1, w_2 \dots w_n$  remain almost  
Static

This is known as "lazy" training  
(Empirical observation)

Writing first order Taylor's  
series around  $w_0 - \phi(x)^T$

$$f(w, x) \sim f(w_0, x) + \underbrace{\nabla_w f(w_0, x)^T}_{(w - w_0) \dots}$$

→ Linear in  $w$  but not linear in  $x$

⇒ Similar to the kernel method

$$\phi(x) \triangleq \nabla_w f(w_0, x)$$

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$



# Neural Tangent Kernel (NTK)

Finding NTK for our Neural Net

$$f_m(w, x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \sigma(a_i^T x)$$

↓  
No. of neurons

$$\nabla_{a_i} f_m(w, x) = \frac{1}{\sqrt{m}} b_i \sigma(a_i^T x) x$$

$$\nabla_{a_i} f(w, x) = \sigma(a_i^T x)$$



$$b_i, \frac{1}{\sqrt{m}}$$

$$K_m(x, x') = K_m^a(x, x') + K_m^b(x, x')$$

↓

NTK

$$K_m(x, x') =$$

$$\sum_{i=1}^m \frac{1}{m} b_i^2 \sigma'(a_i^T x) \sigma'(a_i^T x') (x \cdot x')$$

$$\star \boxed{K_m(x, x') = \frac{1}{m} \sum_{i=1}^m \sigma(a_i^T x) \sigma(a_i^T x')}$$

☆  $\prod_{i=1}^m$

Let  $a_i$ s and  $b_i$ s be independent samples

$$K_m^a(x, x') \xrightarrow{m \rightarrow \infty} K^a(x, x') = \mathbb{E} \left( b^2 \sigma(a^T x) \sigma(a^T x') \right)$$

Expectation

$$K_m^b(x, x') \xrightarrow{m \rightarrow \infty} K^b(x, x') = \mathbb{E} \left( \sigma(a^T x) \sigma(a^T x') \right)$$

NTKs for NN with  
infinite width  
( $m \rightarrow \infty$ )

☆ Gradient flow

$$w(t+1) = w(t) - \eta \nabla_w \mathcal{L}(w)$$

$\eta \rightarrow 0$  { condition for  
gradient flow }

$$\frac{w(t+1) - w(t)}{\eta} = -\nabla_w \mathcal{L}(w(t))$$

$$\eta \rightarrow 0$$

$$\frac{dw(t)}{dt} = -\nabla_w \mathcal{L}(w(t))$$

ODE

Gradient flow dynamics  
equation

$$\text{Let } L(w) = \frac{1}{2} \|\hat{y}(w) - y\|^2$$

$$\nabla_w L(w) = \nabla_w \hat{y}(w) (\hat{y}(w) - y)$$

$$\frac{dw(t)}{dt} = - \nabla_w \hat{y}(w) (\hat{y}(w) - y)$$

$$\frac{d\hat{y}(w(t))}{dt} = \underbrace{\nabla_w \hat{y}(w(t))^\top}_{\text{Chain Rule}} \underbrace{\frac{dw(t)}{dt}}_{\text{Chain Rule}}$$

Chain Rule

$$\frac{d\hat{y}(w(t))}{dt} = -\nabla_w \hat{y}(w(t))^T x$$

$$\nabla_w \hat{y}(w) (\hat{y}(w) - y)$$



$$N \nabla K \equiv K(w_0)$$



$$\frac{d\hat{y}}{dt} \approx -K(w_0) (\hat{y}(w(t)) - y)$$

Let ...

$$\Delta U \quad U = y - y$$

$$\frac{dU}{dt} \approx -K(\omega_0) U$$

$$\Rightarrow \boxed{U(t) = U(0) e^{-K(\omega_0)t}} \quad \star$$