③ RBF kernel (Gaussian kernel)

For any $\alpha > 0$    $\emptyset: R^d \rightarrow R^{\infty}$

$K(x, x') \quad \emptyset(x)^T\emptyset(x') = \exp\left(-\frac{\|x-x'\|_2^2}{2\sigma^2}\right)$

$$\boxed{\text{Deep Network 1}}$$

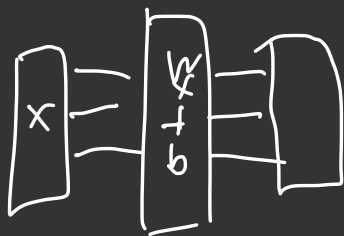$\boxed{\text{Basic structure}}$:

   i'terated linear predictor.

1 layer :   $x \longrightarrow W_1 x + b_1$

2 layers :   $x \longrightarrow W_2 (W_1 x + b_1) + b_2$

3 layers:   $x \longrightarrow W_3 (W_2 (W_1 x + b_1) + b_2) + b_3$

$L$ layers   $x \longrightarrow W_L (\cdots (W_1 x + b_1) \cdots) + b_L$

$$W_2(W_1 x + b_1) + b_2$$
$$= W_2 W_1 x + [W_2 b_1 + b_2]$$

$$W_L \quad (\cdots (W_1 x + b_1) \cdots) + b_L$$
$$= \boxed{(W_L \cdots W_1) x + (b_L + W_L b_{L-1} + \cdots + W_{L \cdots} W_2 b_1)}$$

$$= W^T \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$\downarrow \qquad w \in \mathbb{R}^{d+1}$$

$$W^T_{1:d} = W_L \cdots W_1 \qquad W_{d+1} = b_1 + W_L b_{L-1}$$
$$\qquad\qquad\qquad\qquad\qquad + \cdots + W_L \cdots W_2 b_1$$

$\textcolor{red}{\text{Just a Linear predictor}}$

$\boxed{\text{Activations / Nonlinerities}}$

$$\Pr[Y=1 \mid X=x] = \frac{1}{1 + \exp(-W^T x)} =: \sigma_s(w^T x)$$

$\sigma_s$ logistic or sigmoid function

$$W_2 \sigma(W_1 x + b) + b_2$$

## Classical deep network

$$x \longrightarrow \sigma_L \left( W_L \sigma_{L-1} \left( \cdots \left( W_2 \sigma_1 \left( W_1 x + b_1 \right) + b_2 \right) \cdots \right) + b_L \right)$$

$$\equiv x \longrightarrow (f_L \circ \cdots \circ f_1)(x) \qquad f_i(z) = \sigma_i(W_i z + b_i)$$

**weights** $(W_i)_{i=1}^{L} \in R^{d_i \times d_{i-1}}$

**biases** $(b_i)_{i=1}^{L}$

**activations** $(\sigma_i)_{i=1}^{L} \qquad \sigma_i : R^{d_{i-1}} \longrightarrow R^{d_i}$

$$(W_i, b_i)_{i=1}^{L} \quad \text{params}$$

**choices**

① binarization: $z \longrightarrow \mathbb{1}[z \geq 0] \in \{0, 1\}$

② sigmoid: $\sigma(z) := \dfrac{1}{1 + \exp(-z)}$

③ Hyperbolic tangent: $z \longrightarrow \tanh(z)$
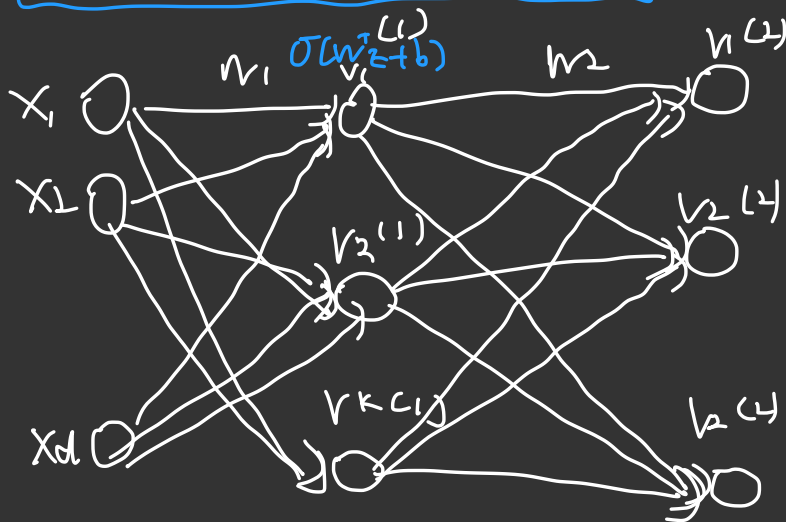
④ Rectified Linear Unit (RELU)
$$\sigma_r(z) = \max\{0, z\} \qquad E.g. \text{ (imagenet)}$$

⑤ identity: $z \longrightarrow z$   last layer when cross entropy (or.

$$\boxed{\text{Multilayer Neural network}} \qquad \left[\boxed{W_1}\;\boxed{\phantom{x}}\right]\left[\boxed{x}\right]^{R}$$



$X_1$    $n_1$   $\sigma(w^T z + b)$    $v_1^{(1)}$    $m_2$    $n^{(2)}$

$X_2$

$v_2^{(1)}$     $v_2 (1)$

$v_K (1)$

$X_d$     $b_2 (1)$

① Columns of $W_1 \in \mathbb{R}^{d \times K}$ : params of original logistic regression models.

② columns of $W_2 \in \mathbb{R}^{K \times F}$: parems of new logistic regression models to combine prediction of original models.

③ $\boxed{\text{multiput}}$ nodes compute $z \to \underline{\sigma(w^T z + b)}$ for some $(W, b)$

④ non-input or non-output units are $\boxed{\text{hidden}}$

# Current "computation graph" perspective.
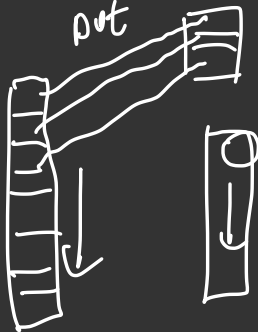
① Edges can pass full tensors



$$\boxed{\sigma(w_1 x + b_1)}$$

② Nodes are more general primitives

③ Edges skip layers

## Convolutional layers    ( < linear )

① 1D convolution
  dot    filter / kernel.
         trainable
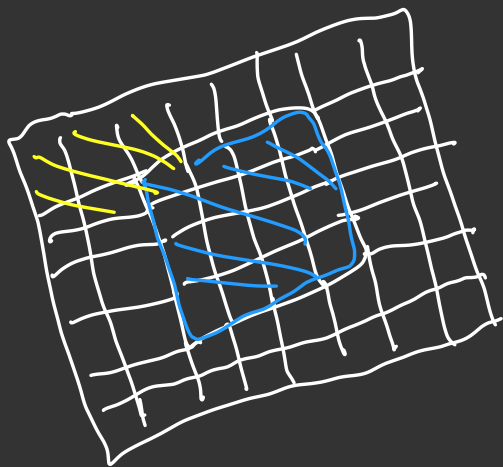


② 2D convolution
  not trainable →  output   input
                   ker

③ padding

kernal
org

④ Strides : step size

channels

extra   below