# 11- learning based vision

the turing test = machine's ability to exhibit intelligent behavior indistinguishable from a human

nearest neighbor classifier → terrible performance for images (does not consider shift, illumination, ..)
↳ training = just remember / keep all the training data $O(n)$ → expensive
↳ predicting = smallest sum of absolute value differences between pixels → return its label
　　　　　　　　　　　　　　　　↑ or square　　　　　　　　　　　　　　　↑ or take k ones

★ CNN's have expensive training, cheap test evaluation

★ do not use test set to determine hyperparameters → use cross validation in train data

challenges in visual recognition = camera pose, illumination, deformation, occlusion, background
clutter, intraclass variations (dog in it)

$$W \cdot x = y \rightarrow f(x,w) + b$$
weights　input　class scores
　　　　　　　　10×1　　10×1

10 × 3072　　32×32×3 = 3072×1

loss functions → hinge loss (tries to find max margin)
↳ multiclass sum loss = $L_i = \sum_{j \neq y_i} max(0, s_j - s_{y_i} + 1)$　$s = f(x_i, w)$ ⟶ $L = \frac{1}{N} \sum_i L_i$

loss of one image = $\sum_{j \neq y_i} max(0, \text{other\_class\_scores} - \text{actual\_class\_score} + 1)$

ex: $\begin{bmatrix} 3.2 \\ 5.1 \\ -1.7 \end{bmatrix}$ actual class $= \underbrace{max(0, 5.1 - 3.2 + 1)}_{2.9} + \underbrace{max(0, -1.7 - 3.2 + 1)}_{0} = 2.9$ (take the average loss for all images)

weight regularization = $L = \frac{1}{N} \sum_i^N L_i + \lambda \cdot R(w) \rightarrow L_1, L_2, L_1 + L_2$ (elastic net), max norm, dropout
　　　　　　　　　　　　　　↳ lambda = regularization strength (hyperparameter)

↳ softmax classifier = normalize log probabilities of classes　$\frac{e^{s_k}}{\sum_j e^{s_j}}$ to maximize → $-log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) = L_i$
　　　　　　　　　　　　　　probabilities [0,1]　　　　　　　　　　　　　　　　　minimize this
ex: $\begin{bmatrix} 3.2 \\ 5.1 \\ -1.7 \end{bmatrix}$ actual class $\xrightarrow[normalize]{exp}$ $\begin{bmatrix} 0.13 \\ 0.87 \\ 0 \end{bmatrix} \rightarrow -log(0.13) = 0.89$
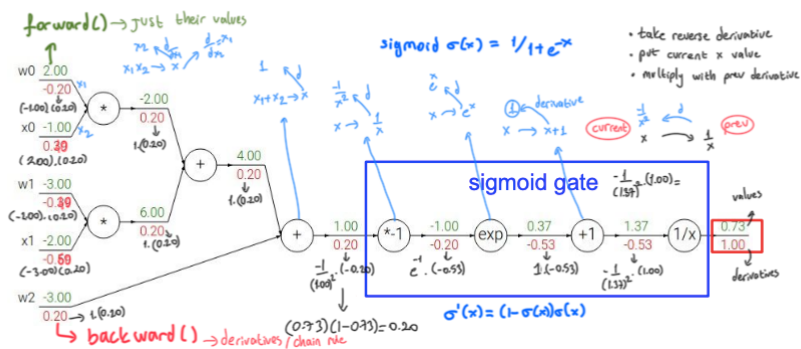
optimization
↳ numerical gradient = computing the gradients for each parameter and observing the change
approximation, slow, easy to write
↳ analytic gradient = exact, fast, but error-prune (because of derivation with math)
↳ gradient check = ensuring the gradients computed by the backpropagation (analytic gradient)
are accurate by numerical gradients
↳ jacobian matrix = matrix of partial derivatives

## activation functions
↳ sigmoid = not centered, vanishes for high and low values, expensive to compute, don't use
↳ tanh = centered at 0 ✓, vanishes, expensive
↳ relu = not centered, does not vanishes for positive values, converges x6 faster, use it ✓

weight initialization = $W = random.randn(input\_size, output\_size) / np.sqrt(input\_size)$
batch normalization = reduces the dependency on initialization    $\hat{x}_i \leftarrow \dfrac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$   $y_i \leftarrow \gamma \hat{x}_i + \beta$

↳ after fully connected layer and before non-linearity           (scale & shift)
↳ during the test time avg means of batches and variances are used

## convolutional neural networks
• 32×32×3 image ✱ 5×5×3 filter ⟶ 28×28×1 activation map
• 32×32×3 image ✱ 6 5×5×3 filters ⟶ 28×28×6 activation map
  ‿‿‿ these must be same
                                            ↳ 10 5×5×6 → 24×24×10
                                               (filter)        (map)

output size = (N−F) / stride + 1
↳ with padding = (2P+N−F) / stride +1

example = input volume = 32×32×3   10 5×5 filters with stride 1, padding 2
↳ output size = (32+2.2−5)/1 +1 = 32 → 32×32×10
↳ number of parameters = 10 ✱ (5×5×3+1) = 10(75+1) = 760
                        ↳ 75 dim dot product in each step

✱ 1×1 convolution can merge or extend the dimensions (can be usefull)

✱ recent trend towards smaller filters, deeper architectures, getting rid of pooling and fully connected layers (just convolution)

upsampling = used to increase size in the next step and in semantic segmentation