

Last time

- convolutional layers
 - pooling layers
 - smaller parameter counts vs dense layers
- } capturing spatial structures

Today

- other important types of layers
 - [→ Batch Normalization
 - Dropout
 - [→ Residual Layer
- } different behavior at train & test time
- can learn to be removed during training

Batch Normalization

→ a "special" layer type which does not follow the linear + nonlinear piece recipe

→ is usually applied either: (1) after the linear piece and before the non-linear piece of another layer, or (2) after the non-linear piece of another layer.

$$(1) \sigma(\text{BatchNorm}(z(x)))$$

$$(2) \text{BatchNorm}(\sigma(z(x)))$$

→ was introduced as a fix to the gradient explosion / gradient vanishing problem, but is perhaps best thought of as a noising / regularizing layer.

→ Idea: center and scale the activations of a layer to prevent the drift of this distribution to extremes during training OR reparameterize the activations of each observation as measuring the relative differences between it and members of the batch.

input: $X = [n, p]$ ^{e.g.} output of previous layer
batch size

Behavior during training:

(1) Obtain the batch mean \bar{X}_B and batch variance S_B^2 feature-wise:

$$\bar{X}_B = \frac{1}{n} \sum_{i=1}^n X_i \in \mathbb{R}^p \quad \triangle! \text{ highly dependent on batch size}$$
$$S_B^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_B)^2$$

(2) Update the moving average estimators of the population mean and variance of activations:

$$\hat{\mu} = \lambda \hat{\mu} + (1-\lambda) \bar{X}_B$$

$$\hat{\sigma}^2 = \lambda \hat{\sigma}^2 + (1-\lambda) S_B^2$$

λ is the momentum hyperparameter ($\sim .99$)

(3) calculate normalized activations:

$$\tilde{X}_i = \left(\frac{X_i - \bar{X}_B}{\sqrt{S_B^2 + \epsilon}} \right)$$

(4) learned scale and shift :

$$z_i = \gamma \tilde{x}_i + \beta$$

where $\gamma \in \mathbb{R}^p$, $\beta \in \mathbb{R}^p$ are trainable parameters, initialized to 1 and 0 respectively

→ allows the network to "undo" the scaling

Behavior during testing :

$$z_i = \gamma \left(\frac{x_i - \hat{\mu}}{\sqrt{\hat{\sigma}^2 + \epsilon}} \right) + \beta$$

→ replace batch mean and variance estimates with estimates of population statistics.

→ why moving average? distribution may change over time during training

Related layers

→ Layer Norm

Dropout

- a "special" layer that does not have the usual linear → nonlinear structure
- is usually applied before a Batch Normalization layer e.g. $\text{BatchNorm}(\text{Dropout}(\sigma(\mathcal{Z}(X))))$
- Idea: Noise a layer's outputs by "dropping" each neuron (i.e. setting its value to zero) for each observation with a probability d .

Behavior during training:

$$\text{Dropout}_d(X) = X \cdot M$$

$[n \times p] \quad [n \times p]$

element-wise multiplication

$$M_{ij} = \mathbb{1} \{ U_{ij} > d \} , \quad U_{ij} \sim \text{uniform}(0, 1)$$

Behavior during testing:

$$\text{Dropout}_d(X) = X \cdot (1 - d)$$

- want the expected magnitudes of the outputs to be the same during training and testing

Behavior during training (v2):

$$\text{Dropout}_d(X) = X \cdot M / (1 - d)$$

Behavior during testing (v2):

$$\text{Dropout}_d(X) = X$$

→ avoids one computation during inference (where latency can be a big issue)

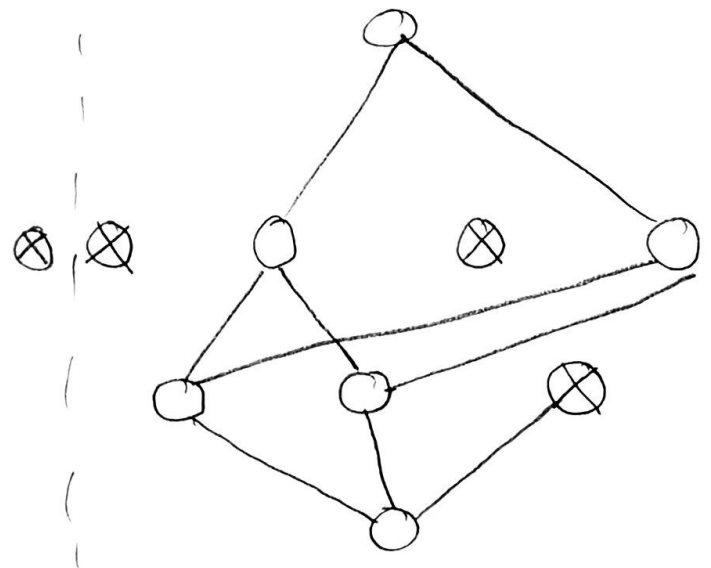
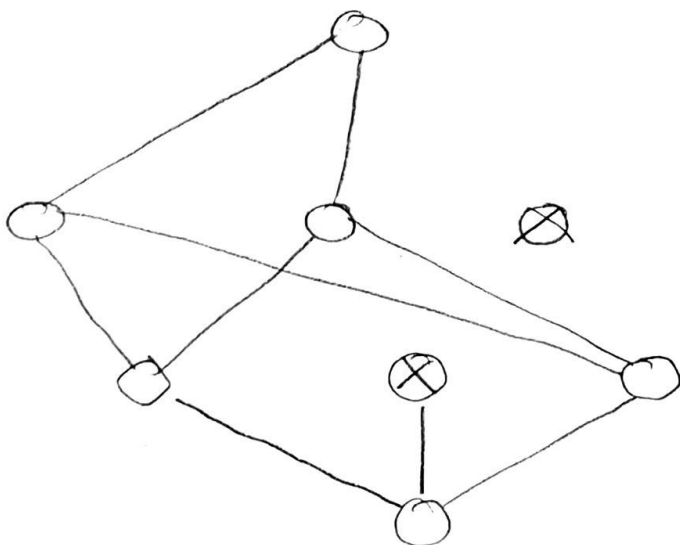
→ Interpretations:

Related: Stochastic Depth (^{layer}drop)
Drop Connect (weight drop)

→ prevent over-reliance on any particular feature

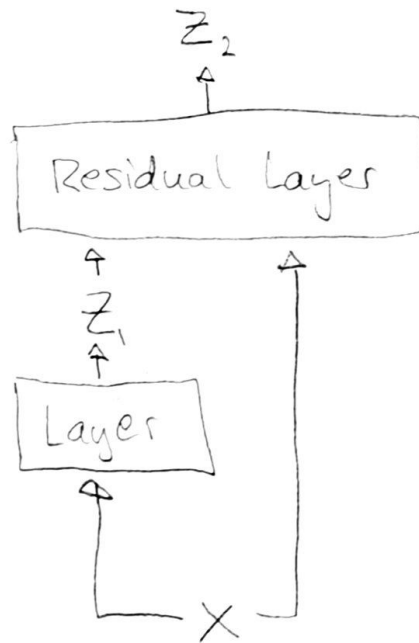
→ encourage utilization of all features

→ learn features that are useful for all sub-networks of the trained model



Residual Layer

→ Idea: Enable a model to choose/learn the optimal number of layers needed during training.



$$Z_2 = Z_1 + X \quad (\text{where } + \text{ is element-wise addition})$$

→ Note: Z_1 and X must have the same shape

→ A key innovation which enabled training of very deep (> 100 layer) networks, overcoming vanishing/exploding gradient problem.