

Deep Learning

8. Regularization. Dropout. Batch normalization. CNN architectures.

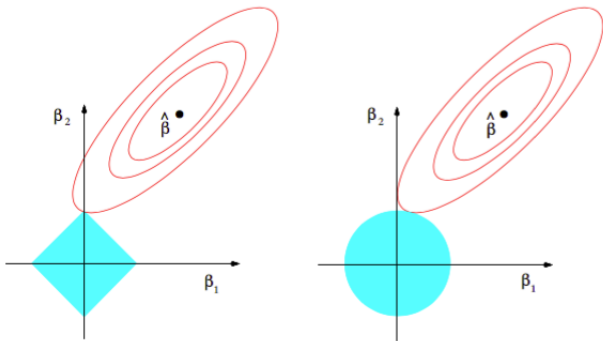
Viacheslav Dudar

Taras Shevchenko National University of Kyiv

2018

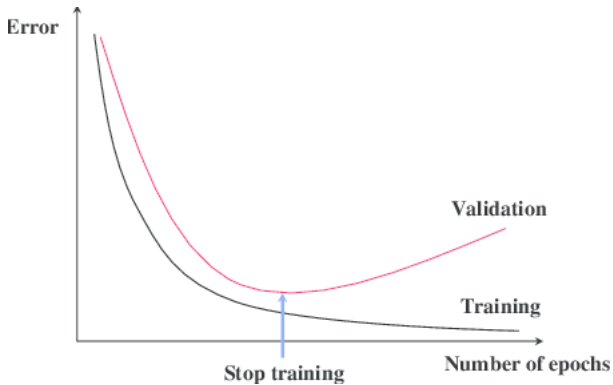
Classical regularization methods

- l_1 , l_2 regularizations



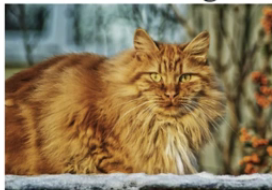
Drawback: no reasonable way to choose regularization coefficient
Partial solution: MacKay's quick and dirty method

Early stopping



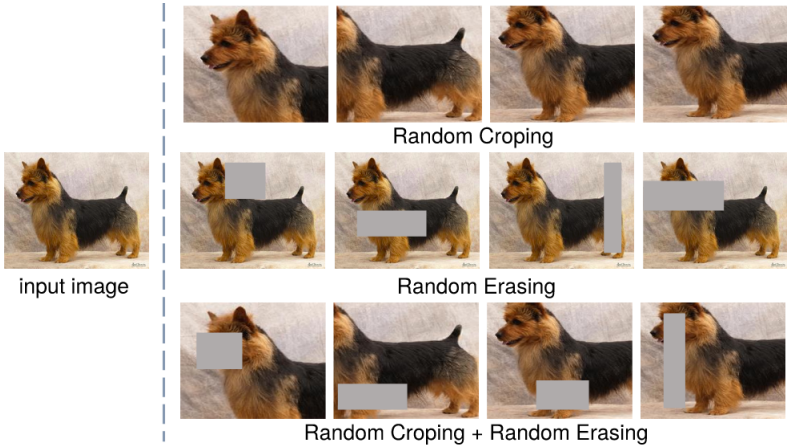
Data augmentation

Mirroring, or horizontal flipping



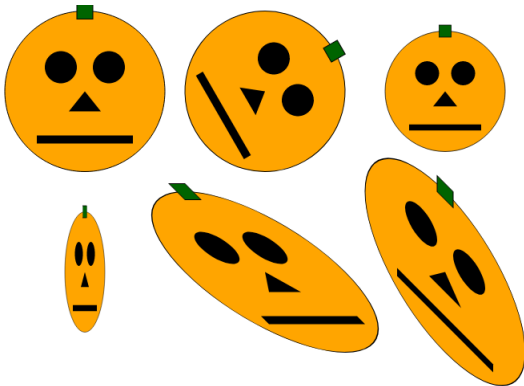
Data augmentation

Random cropping / erasing



Data augmentation

Affine image transformations



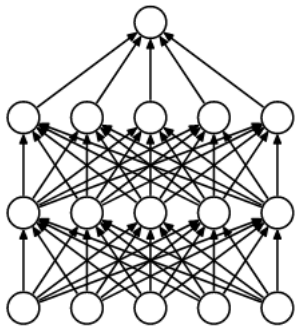
Are not always efficient because they change local statistical properties of pixels that could decrease generalization level.

Data augmentation

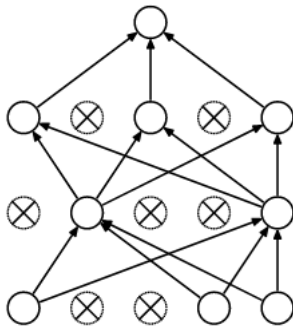
Most popular data augmentation settings:

- Random horizontal flip
- Random image crop with zero padding (image size does not change)
- Random image crop (with image size decrease) at train time; averaging of 5 images (4 corner and 1 center image) at test time

Dropout

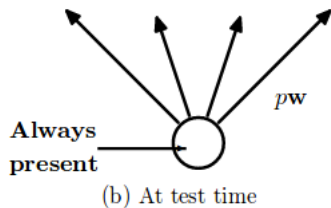
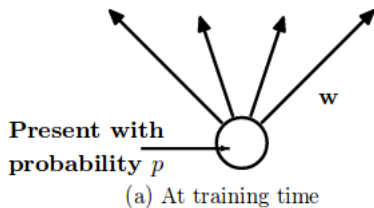


(a) Standard Neural Net



(b) After applying dropout.

Dropout: train-test time



Dropout notes

- Corresponds to averaging of exponential number of smaller models
- Does not have model averaging interpretability if applied to earlier layers, but still improves generalization
- Prevents coadaptation of features
- Popular drop rates are 0.2 - 0.5

Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

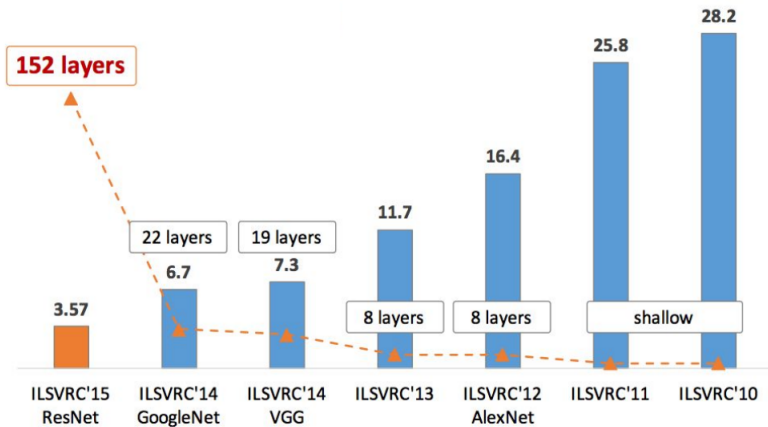
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

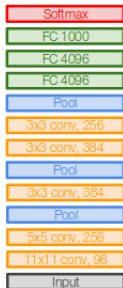
Batch norm properties

- Speeds up training
- Improves generalization
- Is usually used before nonlinearity
- In case of convolutional networks normalization and scale/shift coefficients are shared over spatial dimensions
- Typical sequences of layers in CNNs are Conv-BN-Relu or BN-Relu-Conv

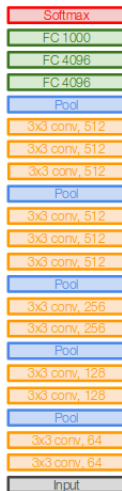
ImageNet competition results



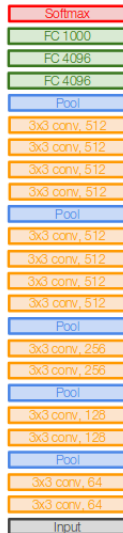
VGG



AlexNet



VGG16

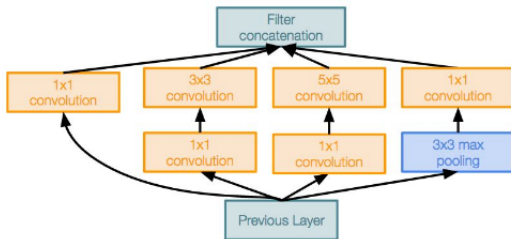


VGG19

VGG - details

- 3*3 filters, padding 1, stride 1
- Dropout + averaging at test time
- 138 M params

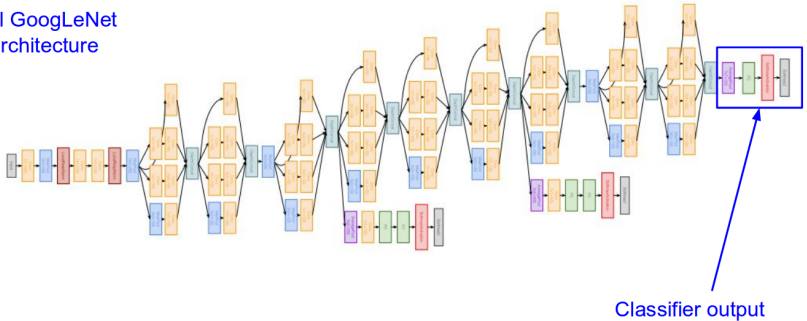
GoogLeNet



Inception module

GoogLeNet

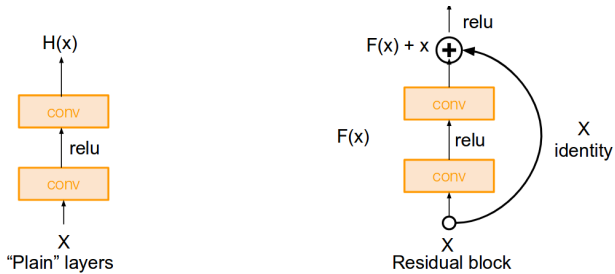
Full GoogLeNet
architecture



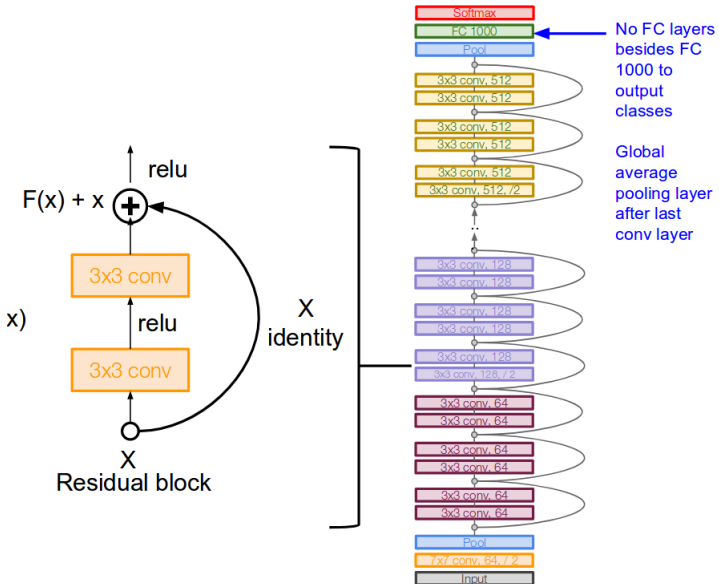
GoogleNet - details

- No fully connected layers at the end
- 5M parameters
- ILSVRC14 classification winner (6.7% top 5 error)
- 22 layers

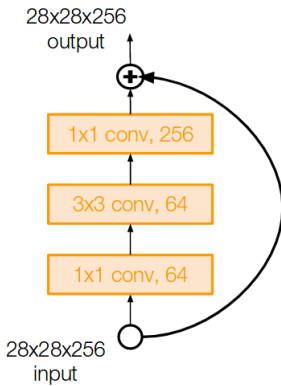
ResNet



ResNet



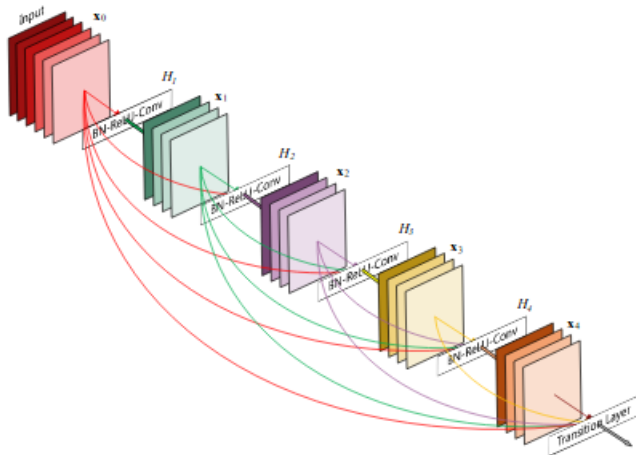
Bottleneck layers



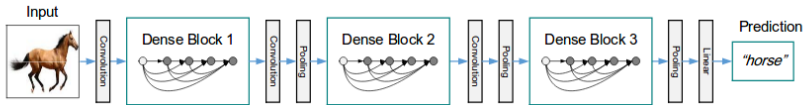
ResNet details

- 34, 50, 101, 152 layers
- Able to train very deep networks without degrading
 - **1st places in all five main tracks**
 - ImageNet Classification: *"Ultra-deep"* (quote Yann) **152-layer** nets
 - ImageNet Detection: **16%** better than 2nd
 - ImageNet Localization: **27%** better than 2nd
 - COCO Detection: **11%** better than 2nd
 - COCO Segmentation: **12%** better than 2nd

Dense net



Dense net



Dense net

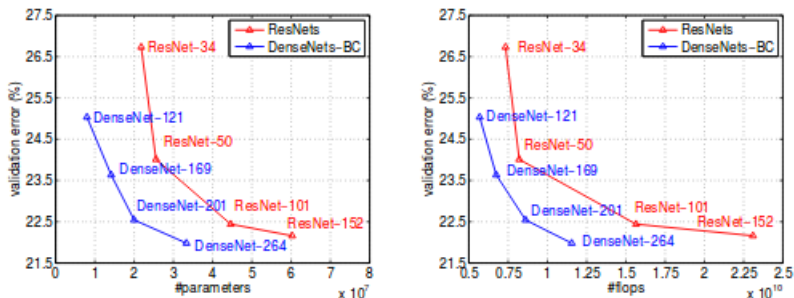


Figure 3: Comparison of the DenseNets and ResNets top-1 error rates (single-crop testing) on the ImageNet validation dataset as a function of learned parameters (*left*) and FLOPs during test-time (*right*).

Dense net

- Dense blocks have such sequences of layers:
BN-Relu-Conv1*1-BN-Relu-Conv3*3-Concat
- Transition blocks: BN-Relu-Conv1*1-Pool with transition rate 0.5
- Gradient freely flows backward from the output to the input
- Model with more layers includes model with less layers
- More layers is associated with better training generalization

Dense net results

Method	Depth	Params	C10	C10+	C100	C100+
Network in Network [22]	-	-	10.41	8.81	35.68	-
All-CNN [32]	-	-	9.08	7.25	-	33.71
Deeply Supervised Net [20]	-	-	9.69	7.97	-	34.57
Highway Network [34]	-	-	-	7.72	-	32.39
FractalNet [17]	21	38.6M	10.18	5.22	35.34	23.30
with Dropout/Drop-path	21	38.6M	7.33	4.60	28.20	23.73
ResNet [11]	110	1.7M	-	6.61	-	-
ResNet (reported by [13])	110	1.7M	13.63	6.41	44.74	27.22
ResNet with Stochastic Depth [13]	110	1.7M	11.66	5.23	37.80	24.58
	1202	10.2M	-	4.91	-	-
Wide ResNet [42]	16	11.0M	-	4.81	-	22.07
	28	36.5M	-	4.17	-	20.50
	16	2.7M	-	-	-	-
ResNet (pre-activation) [12]	164	1.7M	11.26*	5.46	35.58*	24.33
	1001	10.2M	10.56*	4.62	33.47*	22.71
DenseNet ($k = 12$)	40	1.0M	7.00	5.24	27.55	24.42
DenseNet ($k = 12$)	100	7.0M	5.77	4.10	23.79	20.20
DenseNet ($k = 24$)	100	27.2M	5.83	3.74	23.42	19.25
DenseNet-BC ($k = 12$)	100	0.8M	5.92	4.51	24.15	22.27
DenseNet-BC ($k = 24$)	250	15.3M	5.19	3.62	19.64	17.60
DenseNet-BC ($k = 40$)	190	25.6M	-	3.46	-	17.18