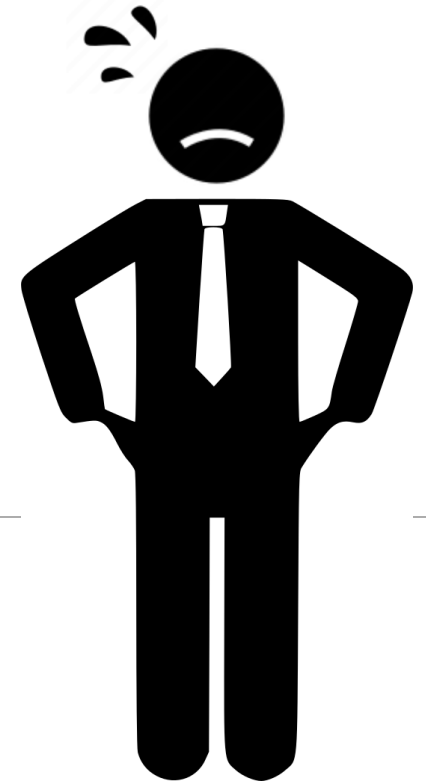


Deep Learning for Small Datasets

DAPHNÉ CHOPARD



Overall Goal (Supervised)

Given X and Y , and assuming there exists a true function $f(\cdot)$ such that

$$Y = f(X) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_\epsilon)$$

the goal is to **estimate a model $\hat{f}(\cdot)$ that approximates $f(\cdot)$** such that

$$\left(Y - \hat{f}(X)\right)^2$$

is minimal.

Overall Goal

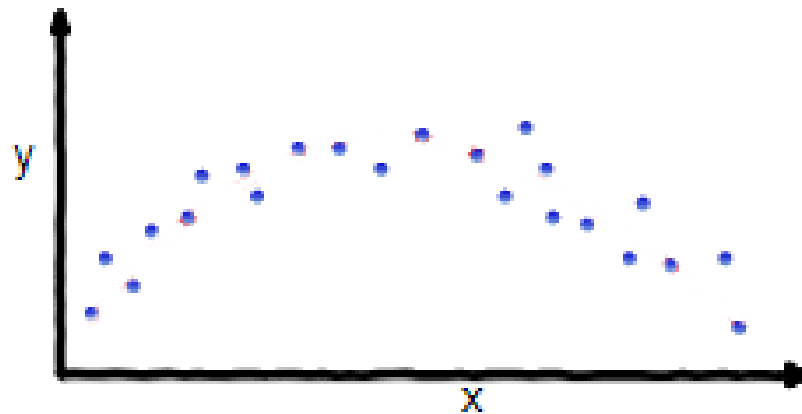
Estimate relationship between input X and output y :

$$\hat{f}(X) = y$$

Overall Goal

Estimate relationship between input X and output y :

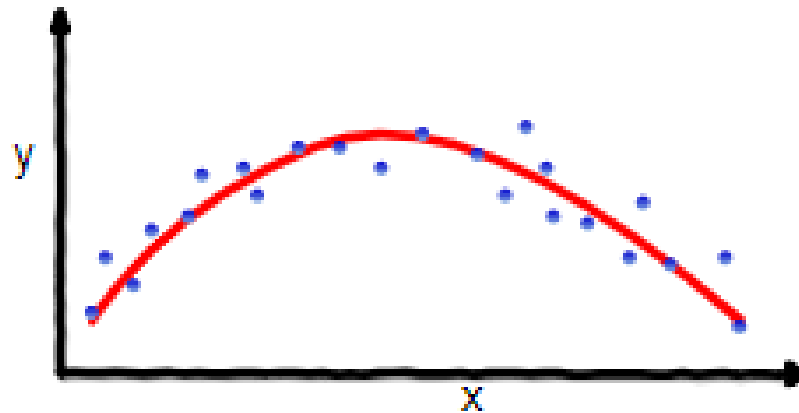
$$\hat{f}(X) = y$$



Overall Goal

Estimate relationship between input X and output y :

$$\hat{f}(X) = y$$



Expected Generalization Error

$$\text{Err}(x) = E \left[\left(Y - \hat{f}(x) \right)^2 \right]$$

Expected Error of Approximation

$$\begin{aligned}\text{Err}(x) &= E \left[\left(Y - \hat{f}(x) \right)^2 \right] \\ &= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2\end{aligned}$$

Expected Error of Approximation

$$\begin{aligned}\text{Err}(x) &= E \left[\left(Y - \hat{f}(x) \right)^2 \right] \\ &= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2 \\ &= f(x)^2 - 2E[\hat{f}(x)]f(x) + E[\hat{f}(x)^2] + \sigma_\epsilon^2\end{aligned}$$

Expected Error of Approximation

$$\text{Err}(x) = E \left[\left(Y - \hat{f}(x) \right)^2 \right]$$

$$= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2$$

$$= E[\hat{f}(x)]^2 + f(x)^2 - 2E[\hat{f}(x)]f(x) + E[\hat{f}(x)^2] - E[\hat{f}(x)]^2 + \sigma_\epsilon^2$$

Expected Error of Approximation

$$\text{Err}(x) = E \left[\left(Y - \hat{f}(x) \right)^2 \right]$$

$$= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2$$

$$= E[\hat{f}(x)]^2 + f(x)^2 - 2E[\hat{f}(x)]f(x) + E[\hat{f}(x)^2] - E[\hat{f}(x)]^2 + \sigma_\epsilon^2$$

$$= \left(E[\hat{f}(x)] - f(x) \right)^2 + E \left[\left(\hat{f}(x) - E[\hat{f}(x)] \right)^2 \right] + \sigma_\epsilon^2$$

Expected Error of Approximation

$$\text{Err}(x) = E \left[\left(Y - \hat{f}(x) \right)^2 \right]$$

$$= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2$$

$$= E[\hat{f}(x)]^2 + f(x)^2 - 2E[\hat{f}(x)]f(x) + E[\hat{f}(x)^2] - E[\hat{f}(x)]^2 + \sigma_\epsilon^2$$

$$= \underbrace{\left(E[\hat{f}(x)] - f(x) \right)^2}_{\text{Squared Bias Error}} + E \left[\left(\hat{f}(x) - E[\hat{f}(x)] \right)^2 \right] + \sigma_\epsilon^2$$

Squared Bias Error

Expected Error of Approximation

$$\text{Err}(x) = E \left[\left(Y - \hat{f}(x) \right)^2 \right]$$

$$= f(x)^2 + E[\hat{f}(x)^2] - 2E[\hat{f}(x)]f(x) + \sigma_\epsilon^2$$

$$= E[\hat{f}(x)]^2 + f(x)^2 - 2E[\hat{f}(x)]f(x) + E[\hat{f}(x)^2] - E[\hat{f}(x)]^2 + \sigma_\epsilon^2$$

$$= \underbrace{\left(E[\hat{f}(x)] - f(x) \right)^2}_{\text{Squared Bias Error}} + \underbrace{E \left[\left(\hat{f}(x) - E[\hat{f}(x)] \right)^2 \right]}_{\text{Variance Error}} + \sigma_\epsilon^2$$

Squared Bias Error

Variance Error

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

High bias ?

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

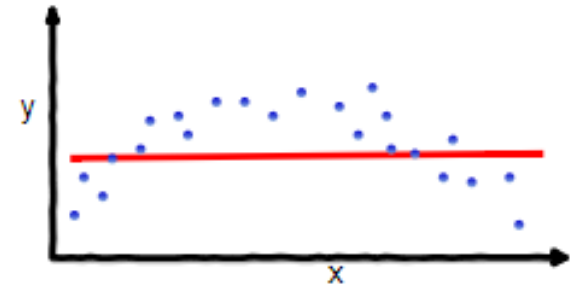
High bias \Leftrightarrow model too simple to fit well

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

High bias \Leftrightarrow model too simple to fit well



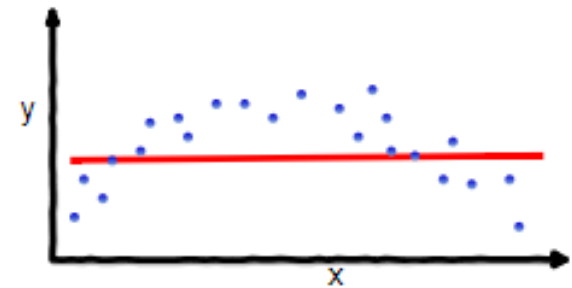
Underfitting

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

High bias \Leftrightarrow model too simple to fit well
 \Leftrightarrow training error large



Underfitting

Bias Error

- How far in general are the predictions from the correct value?

$$\text{Bias}(\hat{f}(x)) = E[\hat{f}(x)] - f(x)$$

High bias	↔	model too simple to fit well
	↔	training error large
Low bias	↔	model complex enough to fit well
	↔	training error low

Variance Error

- How much do the predictions vary between different model realizations?

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

Variance Error

- How much do the predictions vary between different model realizations?

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

High variance?

Variance Error

- How much do the predictions vary between different model realizations?

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

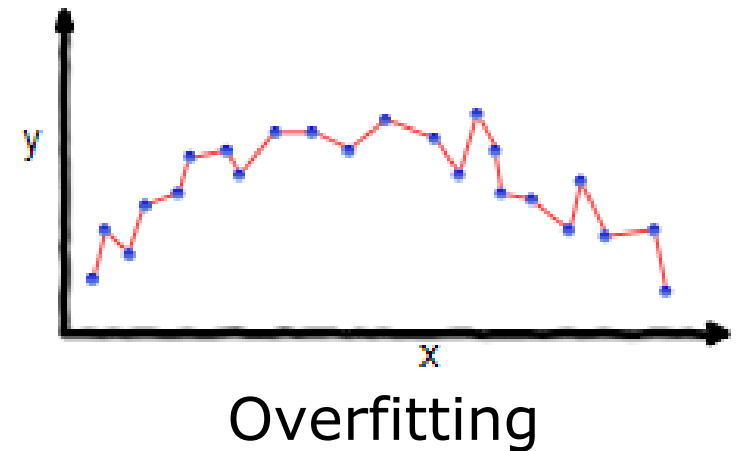
High variance \Leftrightarrow model too complex

Variance Error

- How much do the predictions vary between different model realizations?

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

High variance \Leftrightarrow model too complex

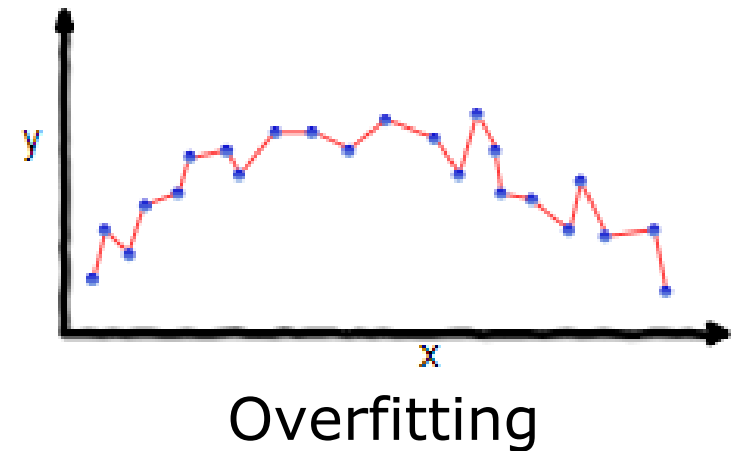


Variance Error

- How much do the predictions vary between different model realizations?

$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

High variance \Leftrightarrow model too complex
 \Leftrightarrow validation error large



Variance Error

- How much do the predictions vary between different model realizations?

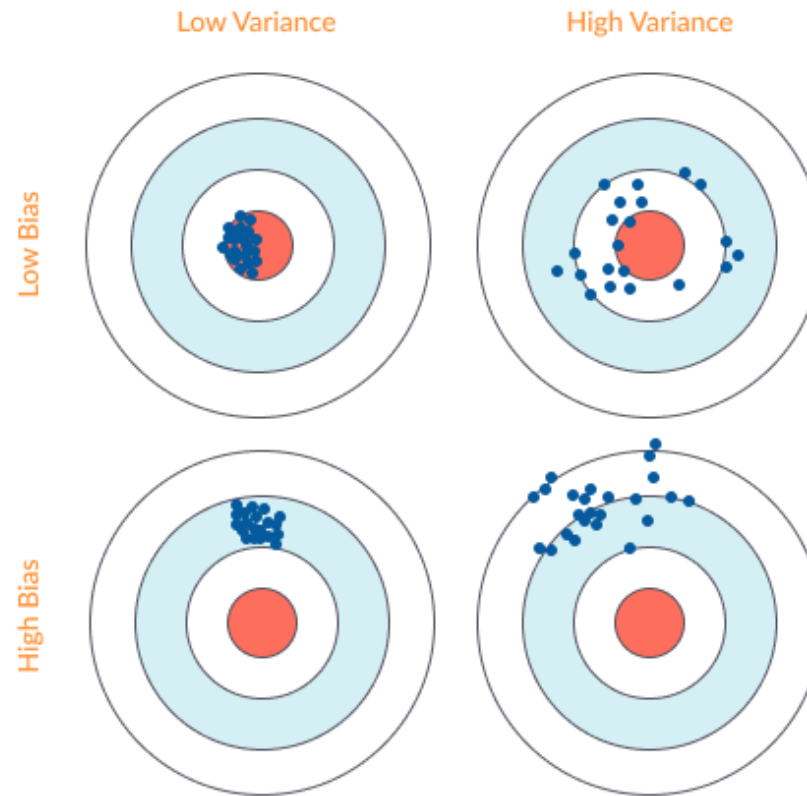
$$\text{Var}(\hat{f}(x)) = E[\hat{f}(x)^2] - E[\hat{f}(x)]^2$$

High variance \Leftrightarrow model too complex to generalize well
 \Leftrightarrow validation error large

Low variance \Leftrightarrow model simple enough to generalize well

Bias-Variance

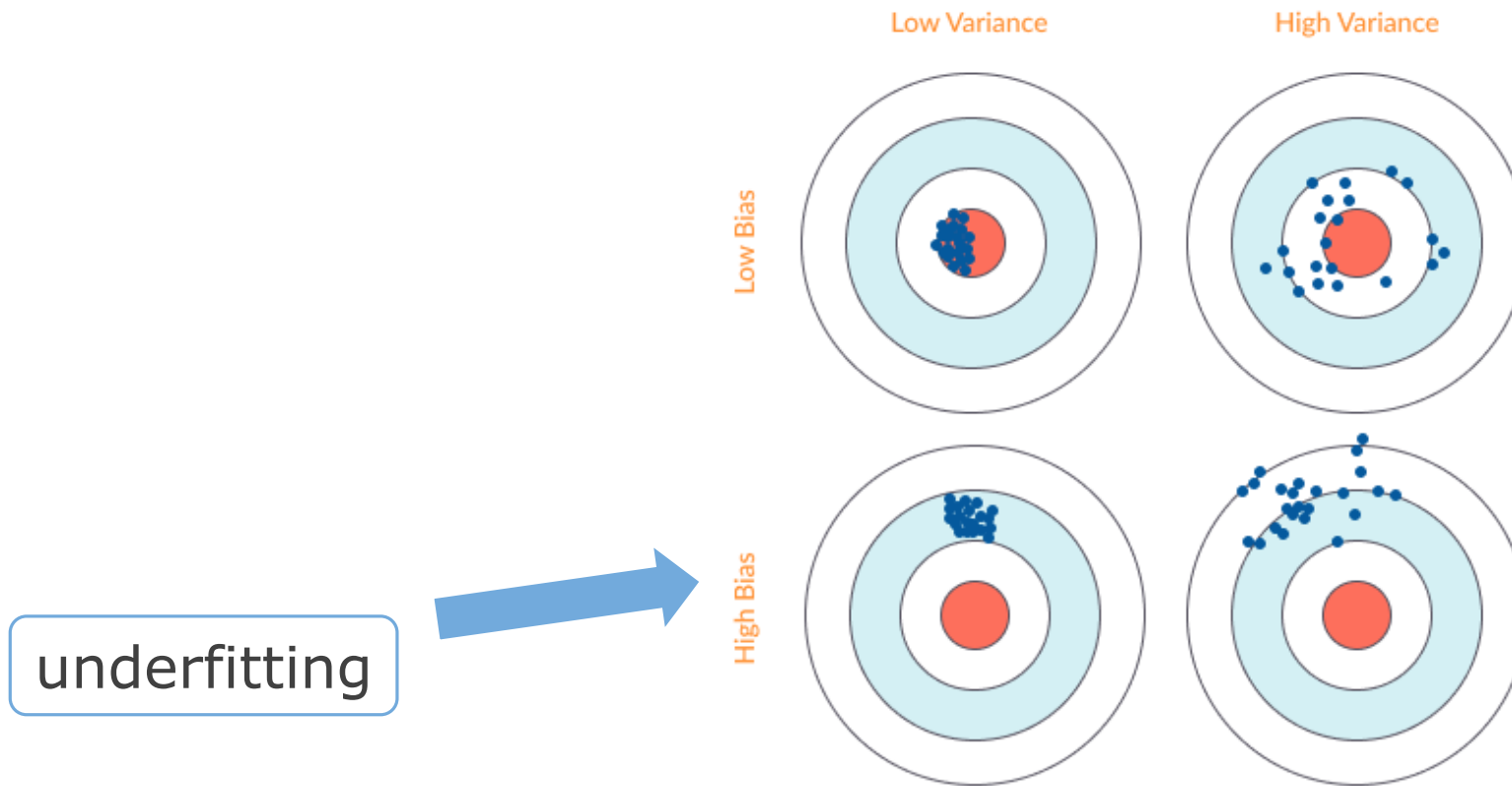
Repeat the model building process multiple times



The centre of the target is a model that perfectly predicts the correct values

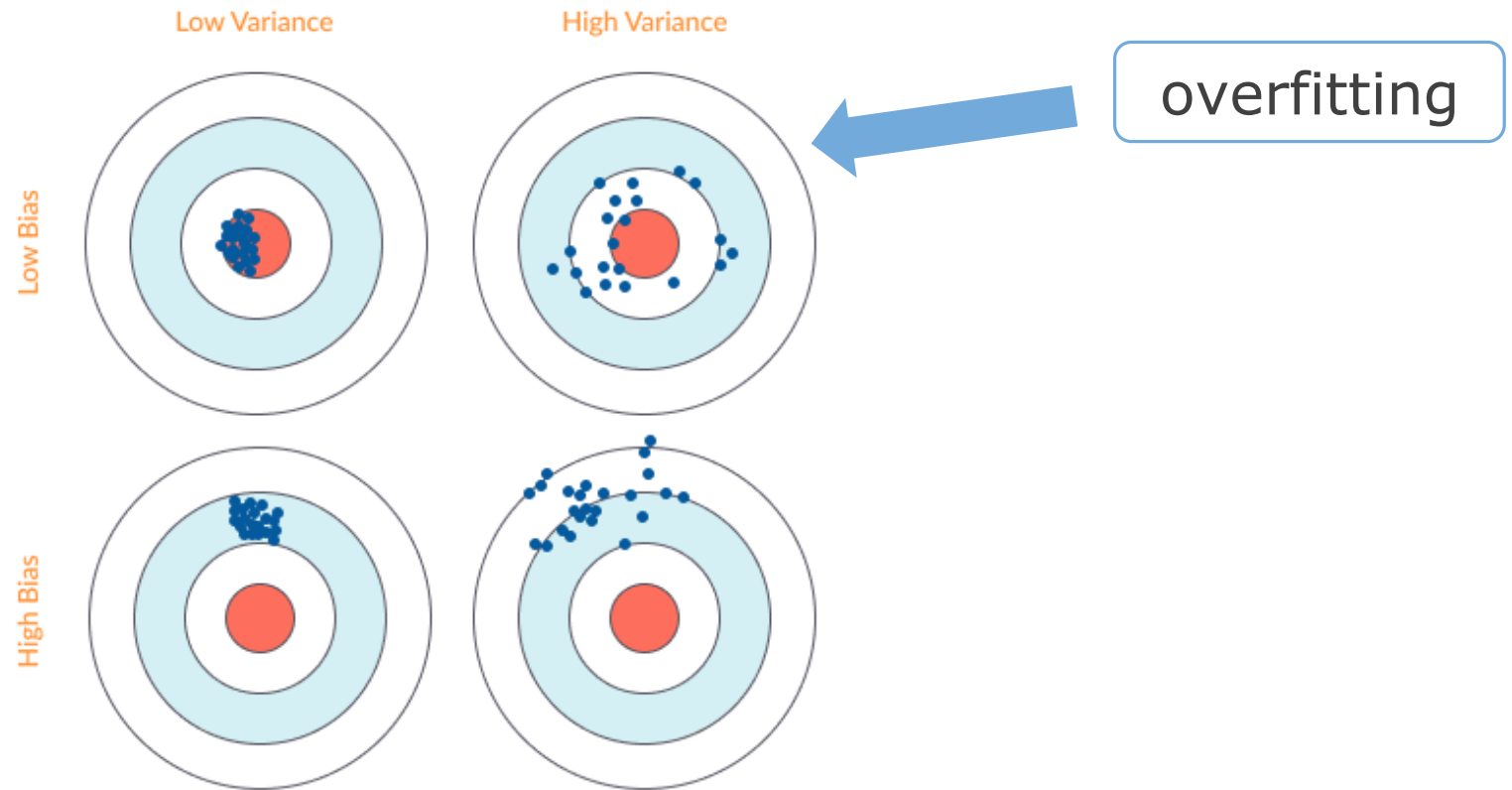
Bias-Variance

Repeat the model building process multiple times

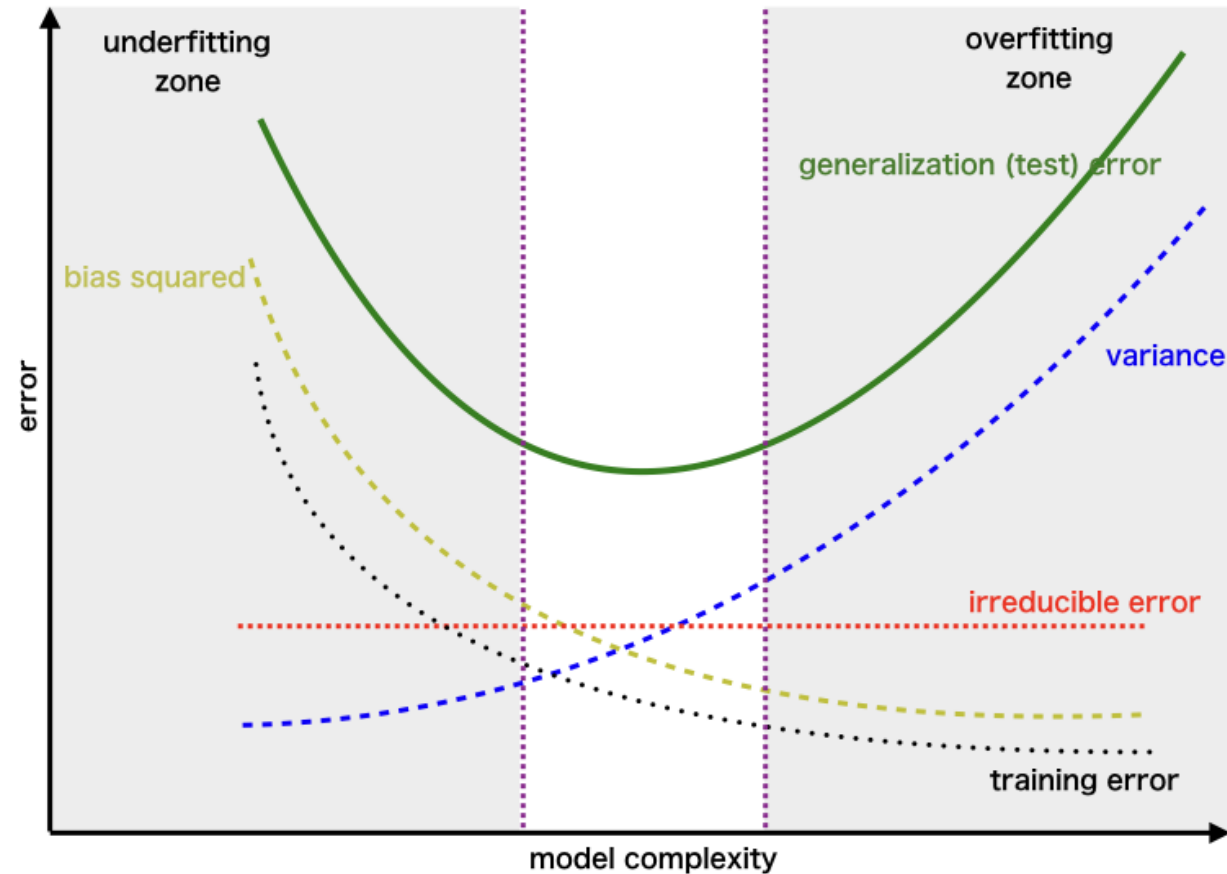


Bias-Variance

Repeat the model building process multiple times






Bias-Variance Trade-Off



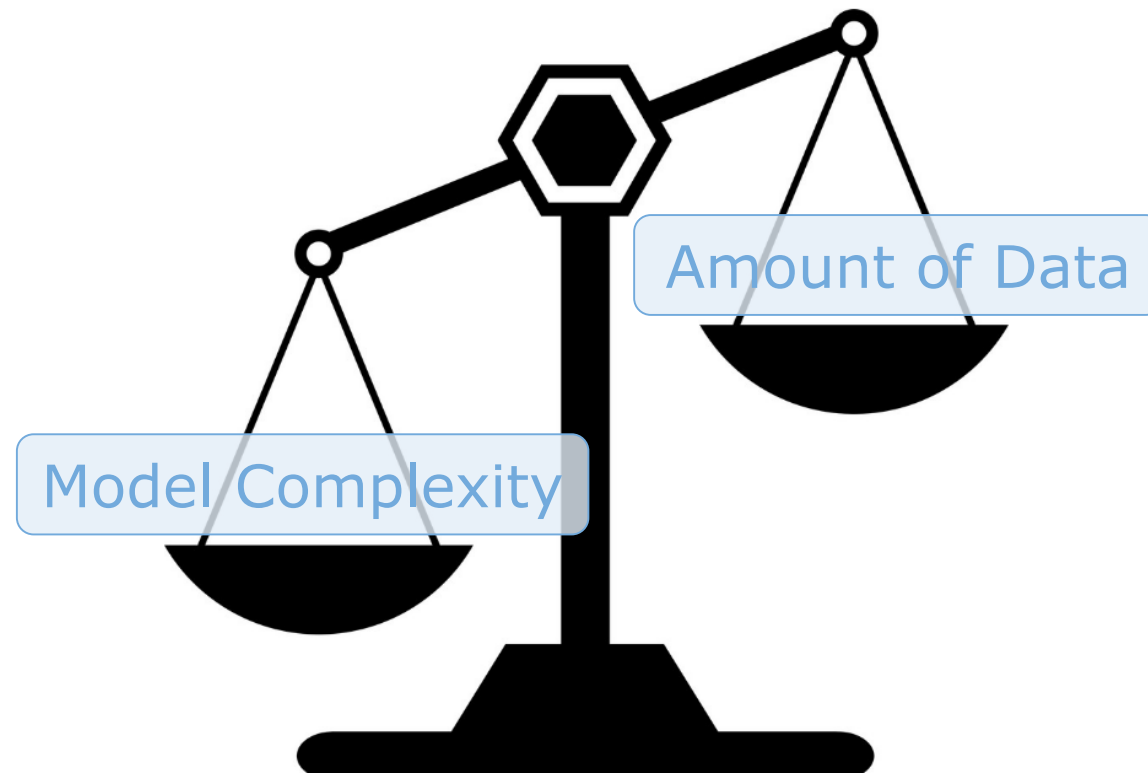
Issues

- Deep neural networks are prone to overfitting (highly complex models)
- Model complexity tied to task complexity

Examples of Competitive DNN

	VGGNet	DeepVideo	GNMT
Used For	Identifying Image Category	Identifying Video Category	Translation
Input	Image 	Video 	English Text 
Output	1000 Categories	47 Categories	French Text
Parameters	140M	~100M	380M
Data Size	1.2M Images with assigned Category	1.1M Videos with assigned Category	6M Sentence Pairs, 340M Words
Dataset	ILSVRC-2012	Sports-1M	WMT'14
	2014	2014	2016

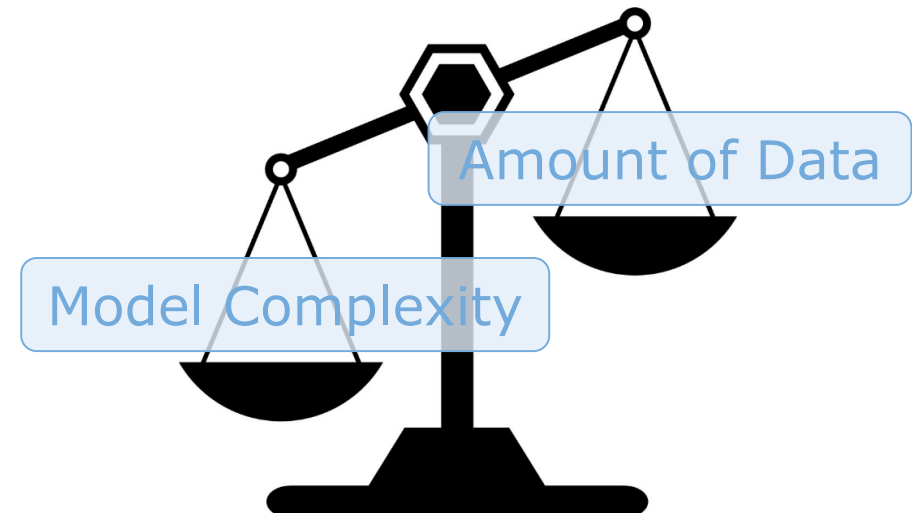
Need to find balance



What if limited data?

- Expensive, labor-intensive to collect
- Usage restriction
 - Sensitive data (confidentiality issues)
- Class imbalance
 - More healthy people than people with a given disease

What if limited data?



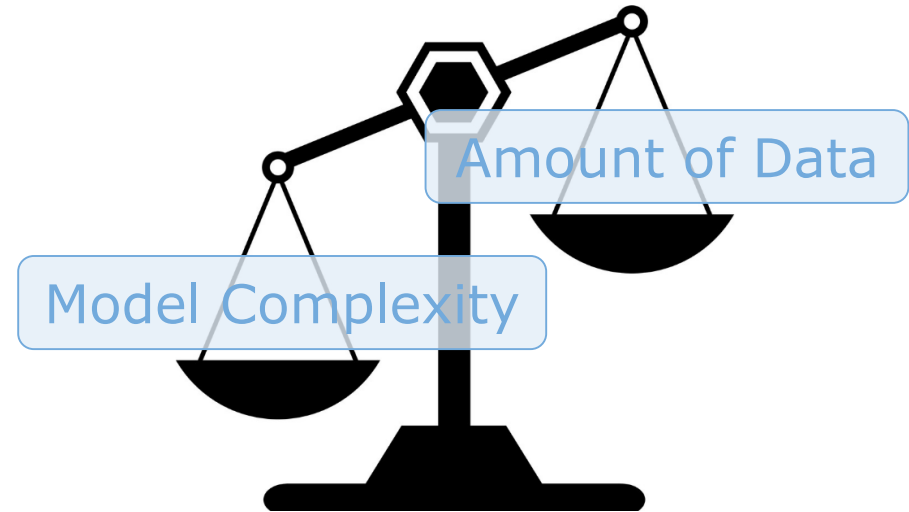
What if limited data?

Reduce Model Complexity

L2 regularisation

L1 regularisation

Dropout



L2 regularisation (weight decay)

- Idea: constrain the network weights by adding a regularization term to the loss function $J(\mathbf{w})$

L2 regularisation (weight decay)

- Idea: constrain the network weights by adding a regularization term to the loss function $J(\mathbf{w})$

$$\underbrace{\tilde{J}(\mathbf{w})}_{\text{New loss function}} = J(\mathbf{w}) + \underbrace{\alpha \|\mathbf{w}\|_2^2}_{\text{Regularization term}}$$

L2 regularisation: gradient

$$\tilde{J}(\mathbf{w}) = J(\mathbf{w}) + \alpha \|\mathbf{w}\|_2^2$$

➤ Gradient update:

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$\mathbf{w}_{\text{new}} = (1 - \eta\alpha) \mathbf{w}_{\text{old}} - \eta \nabla_{\mathbf{w}} J(\mathbf{w}_{\text{old}})$$

L1 regularisation (LASSO)

- Idea: constrain the network weights by adding a regularization term to the loss function $J(\mathbf{w})$

L1 regularisation (LASSO)

- Idea: constrain the network weights by adding a regularization term to the loss function $J(\mathbf{w})$

$$\underbrace{\tilde{J}(\mathbf{w})}_{\text{New loss function}} = J(\mathbf{w}) + \underbrace{\alpha \|\mathbf{w}\|_1}_{\text{Regularization term}}$$

L1 regularisation: gradient

$$\tilde{J}(\mathbf{w}) = J(\mathbf{w}) + \alpha \|\mathbf{w}\|_1$$

➤ Gradient:

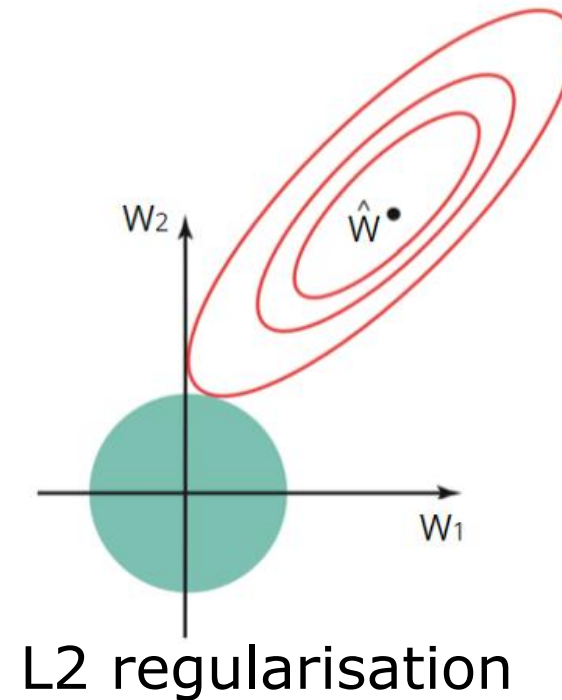
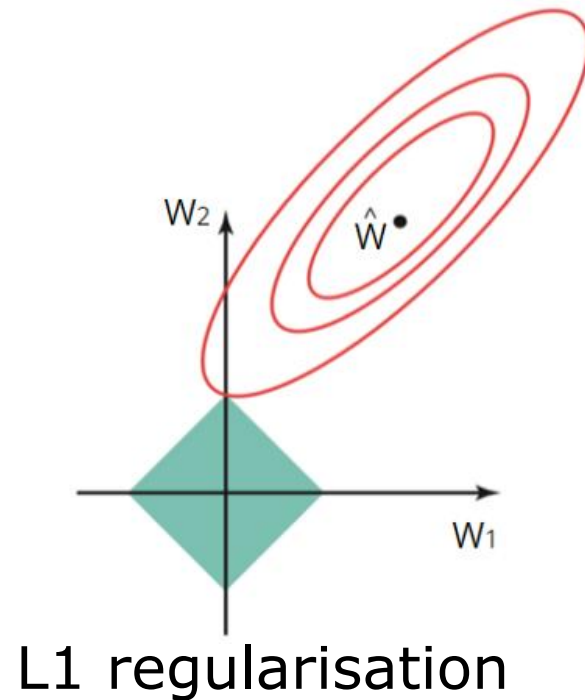
$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}) = \alpha \operatorname{sign}(\mathbf{w}) + \nabla_{\mathbf{w}} J(\mathbf{w})$$

How does regularization work?

- Main idea: smaller weights reduce the impact of the hidden neurons, they become neglectable and the overall complexity of the neural network gets reduced.

How does regularization work?

2D example

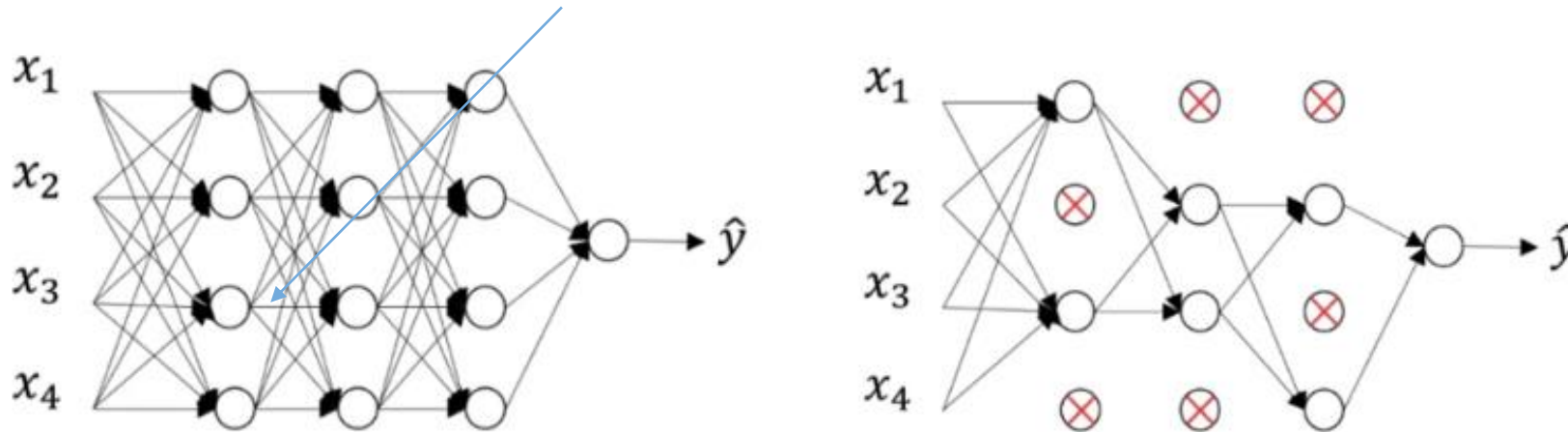


Contours
of loss
function

Constraint
function

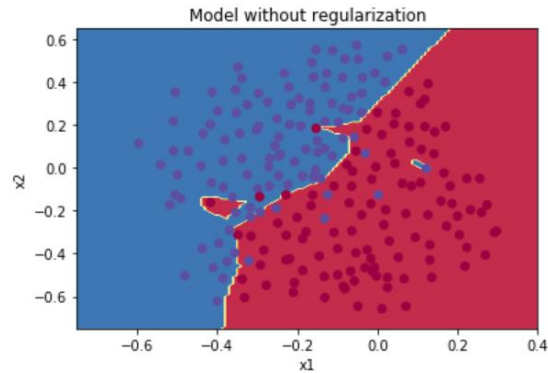
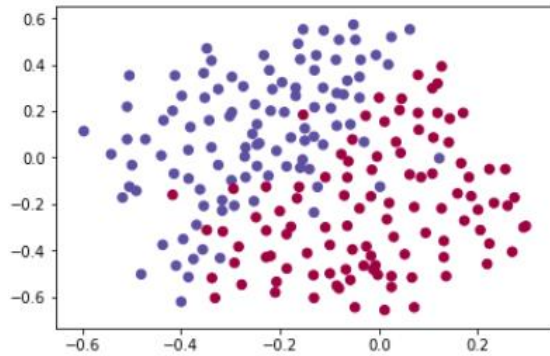
Dropout

- During training turn off a neuron with some probability p

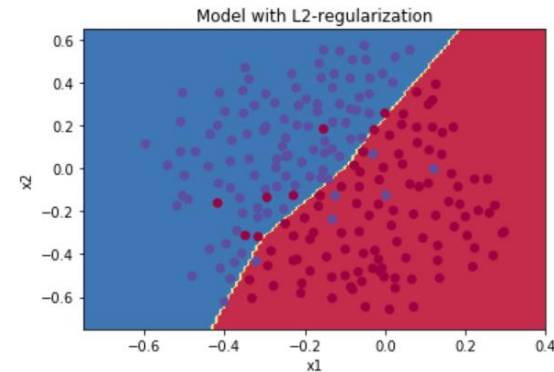


- Idea: The NN will be reluctant to give high weights to certain features, because they might disappear → weights spread across all features making them smaller

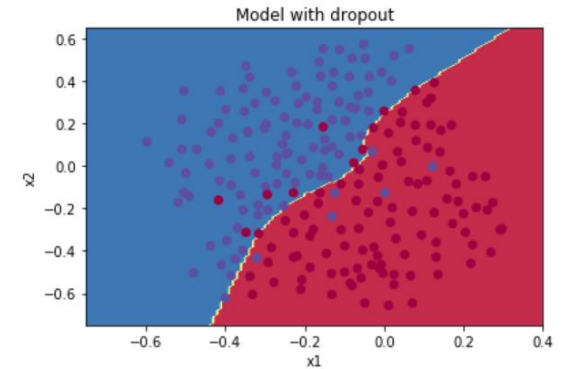
Regularization Example



On the training set:
Accuracy: 0.94786729
On the test set:
Accuracy: 0.915



On the train set:
Accuracy: 0.938388
On the test set:
Accuracy: 0.93

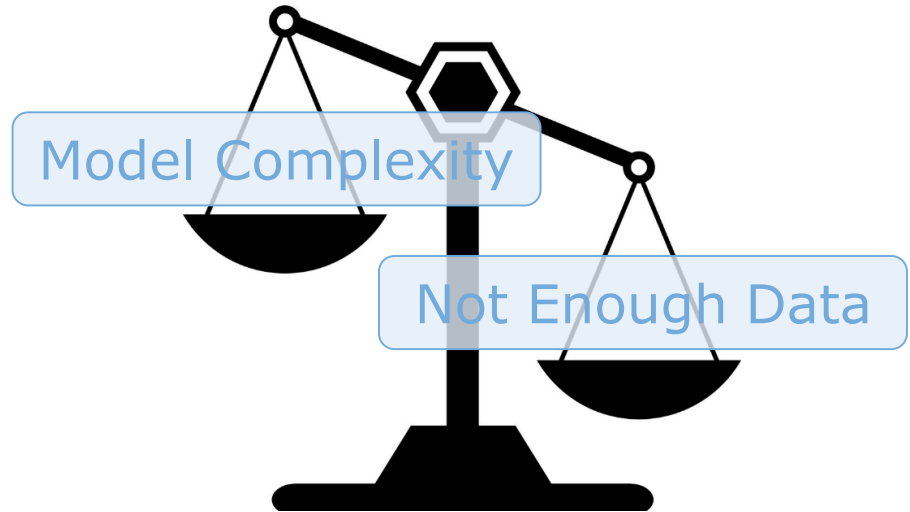


On the train set:
Accuracy: 0.928909
On the test set:
Accuracy: 0.95

What if limited data?

Increase Amount of Data

Data Augmentation



Data Augmentation

- Idea: generate synthetic data from the training data



- The new data must *preserve* the label or the label must be modified accordingly

Data Augmentation Overall

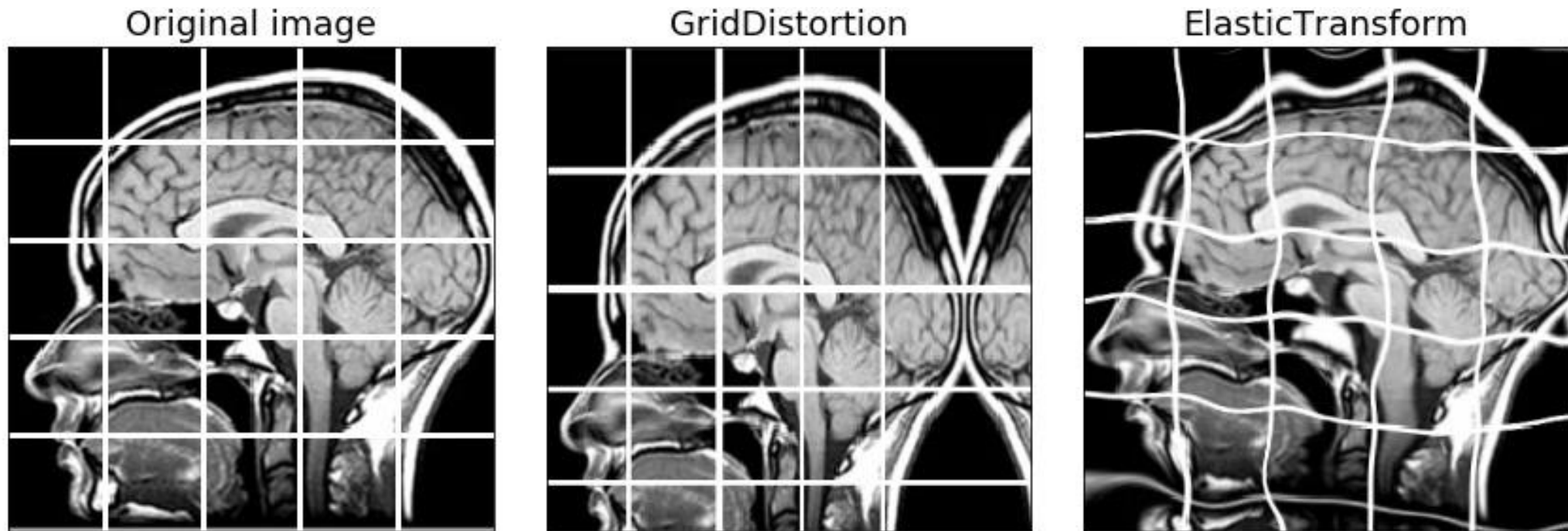
- Increase size and diversity of training data
- Learn invariance to some transformations
- Implicit regularisation effects
- Noising \Leftrightarrow data augmentation

Data Augmentation for Images



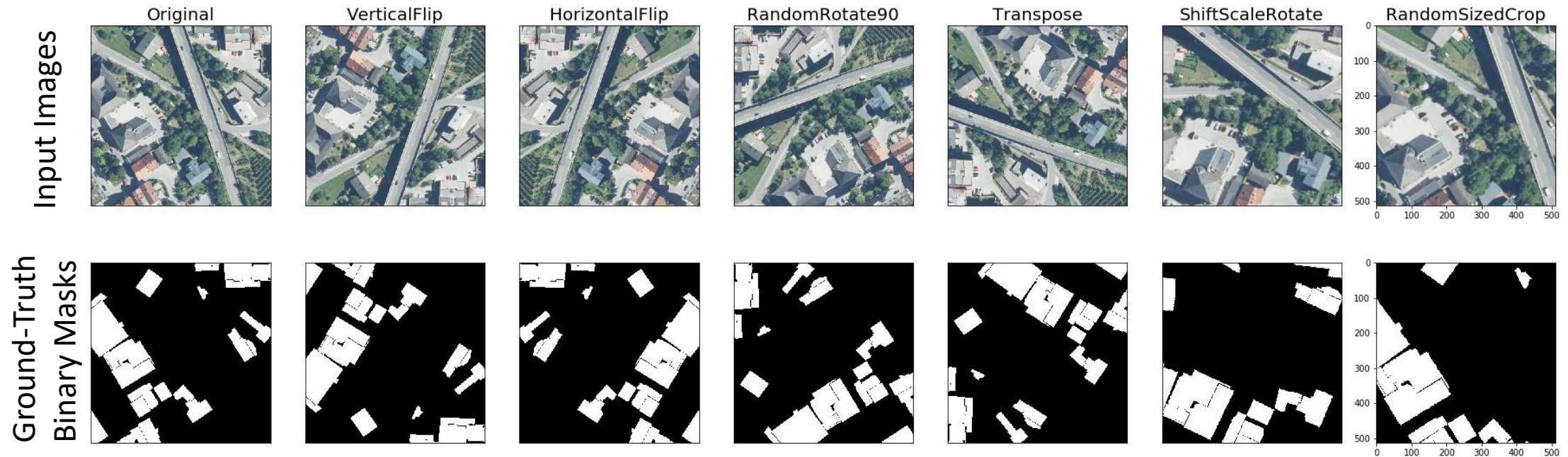
Data Augmentation for Images

- Example of grid transformations commonly used in biomedical image analysis



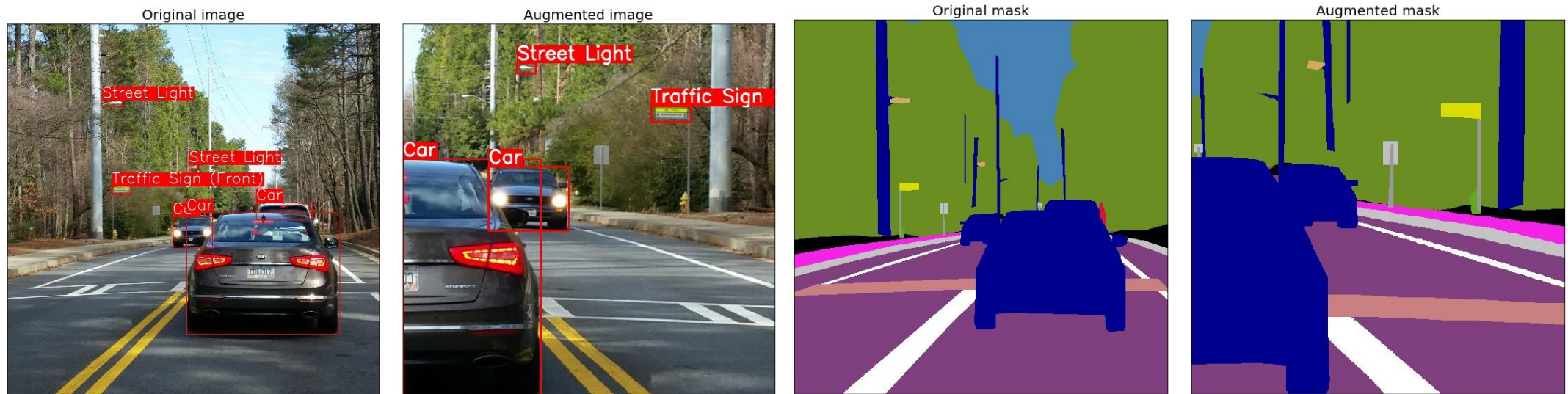
Data Augmentation for Images

- Example of geometry-preserving transforms in a segmentation task



Data Augmentation for Images

- Multiple targets task: an example of applying a combination of transformations to the original image, bounding boxes, and ground truth masks for instance segmentation



Data Augmentation for Images

- Example of results for image classification

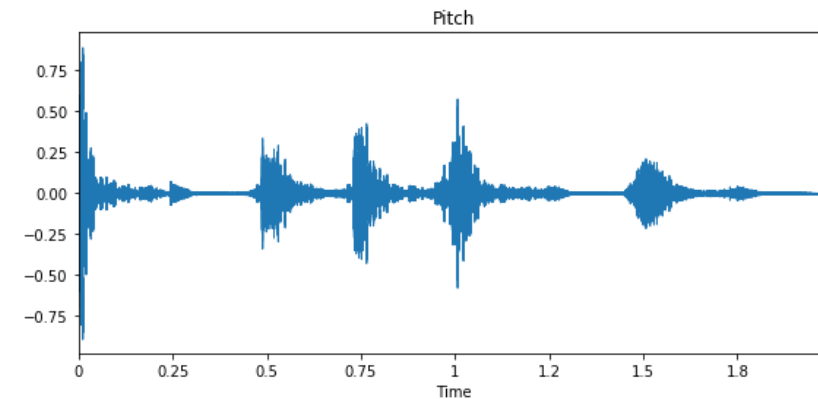
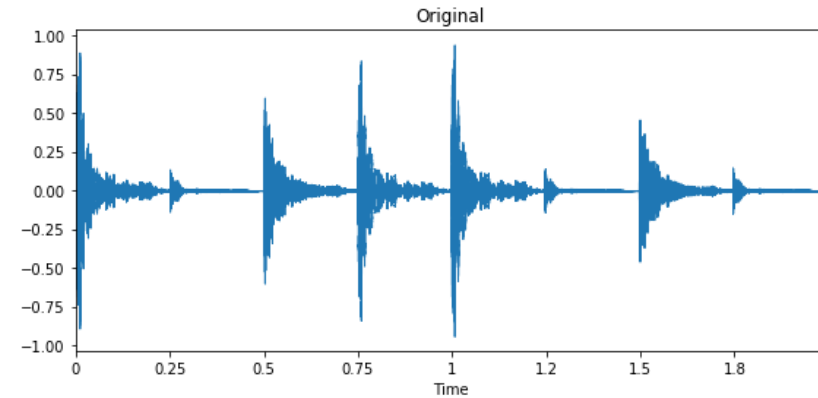


Method	C10	C10+	C100	C100+
ResNet18 [5]	10.63 ± 0.26	4.72 ± 0.21	36.68 ± 0.57	22.46 ± 0.31
ResNet18 + cutout	9.31 ± 0.18	3.99 ± 0.13	34.98 ± 0.29	21.96 ± 0.24
WideResNet [22]	6.97 ± 0.22	3.87 ± 0.08	26.06 ± 0.22	18.8 ± 0.08
WideResNet + cutout	5.54 ± 0.08	3.08 ± 0.16	23.94 ± 0.15	18.41 ± 0.27
Shake-shake regularization [4]	-	2.86	-	15.85
Shake-shake regularization + cutout	-	2.56 ± 0.07	-	15.20 ± 0.21

- Current works focus on *automatically* learning augmentation *schedules*

Data Augmentation for Audio

- Noise Injection
- Time shifting
- Pitch change
- Speed change
- Background noise



Data Augmentation for NLP

“Everyone was in a good mood after enjoying delicious pizzas.”

Data Augmentation for NLP

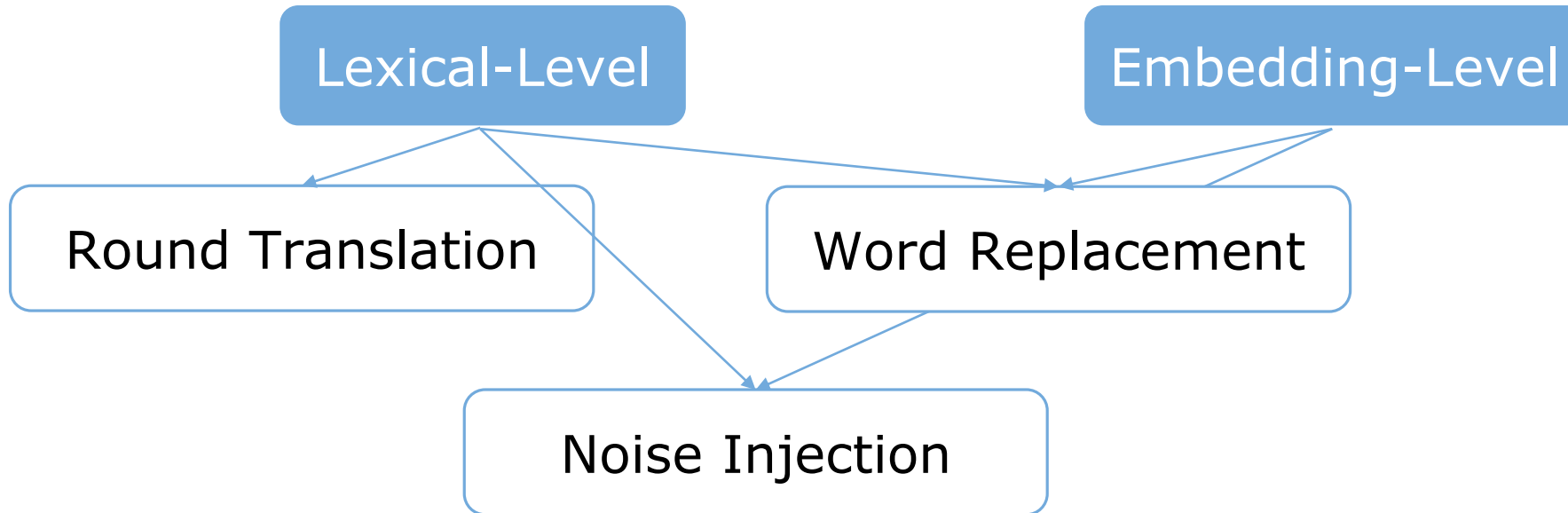
“Everyone was in a good mood after enjoying delicious pizzas.”

Lexical-Level

Embedding-Level

Data Augmentation for NLP

“Everyone was in a good mood after enjoying delicious pizzas.”



Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

“Everyone was in a good mood after enjoying delicious pizzas.”

Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words



“Everyone was in a good mood after enjoying task delicious pizzas.”

Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

“Everyone was in a good ___ after enjoying task delicious pizzas.”



Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words



“Everyone good in a was ____ after enjoying talk delicious pizzas.”

The diagram illustrates the 'Swapping' operation at the lexical level. A curved double-headed arrow is positioned above the words 'good' and 'in' in the sentence below, indicating that these two words are being swapped. The sentence is displayed within a light blue rounded rectangular box.

Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

“Everyone good in a was after enjoying talk delicious pizzas.”

- *Embedding-level:* Adding (e.g., Gaussian) noise to the embeddings

Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

“Everyone good in a was after enjoying talk delicious pizzas.”

- *Embedding-level:* Adding (e.g., Gaussian) noise to the embeddings



Noise Injection

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Inserting, Deleting, Swapping random words

“Everyone good in a was after enjoying talk delicious pizzas.”

- *Embedding-level:* Adding (e.g., Gaussian) noise to the embeddings



Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a good mood after enjoying delicious pizzas.”

Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a cheerful mood after enjoying delicious pizzas.”

Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a good mood after enjoying delicious food.”



Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a cheerful mood after enjoying delicious food.”

- *Embedding-level:* Replace with nearest word embeddings



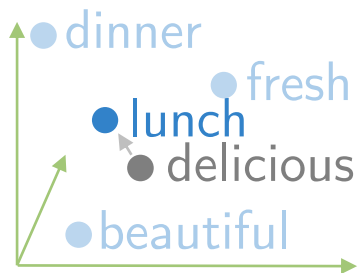
Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a cheerful mood after enjoying delicious food.”

- *Embedding-level:* Replace with nearest word embeddings



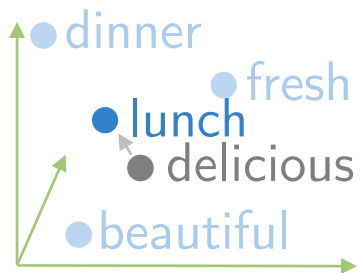
Word Replacement

“Everyone was in a good mood after enjoying delicious pizzas.”

- *Lexical-level:* Replace with synonym, hypernym, language model, ...

“Everyone was in a cheerful mood after enjoying delicious food.”

- *Embedding-level:* Replace with nearest word embeddings



“Everyone was in a good mood after enjoying **lunch** pizzas.”

Round-Translation

- Translation to a target language and then back to source language

“Everyone was in a good mood after enjoying delicious pizzas.”

Round-Translation

- Translation to a target language and then back to source language

“Everyone was in a good mood after enjoying delicious pizzas.”



En → Fr

«*Tout le monde était de bonne humeur après avoir dégusté de délicieuses pizzas.*»

Round-Translation

- Translation to a target language and then back to source language

“Everyone was in a good mood after enjoying delicious pizzas.”



«*Tout le monde était de bonne humeur après avoir dégusté de délicieuses pizzas.*»



“Everyone was in a good mood after tasting delicious pizzas.”

Data Augmentation for NLP

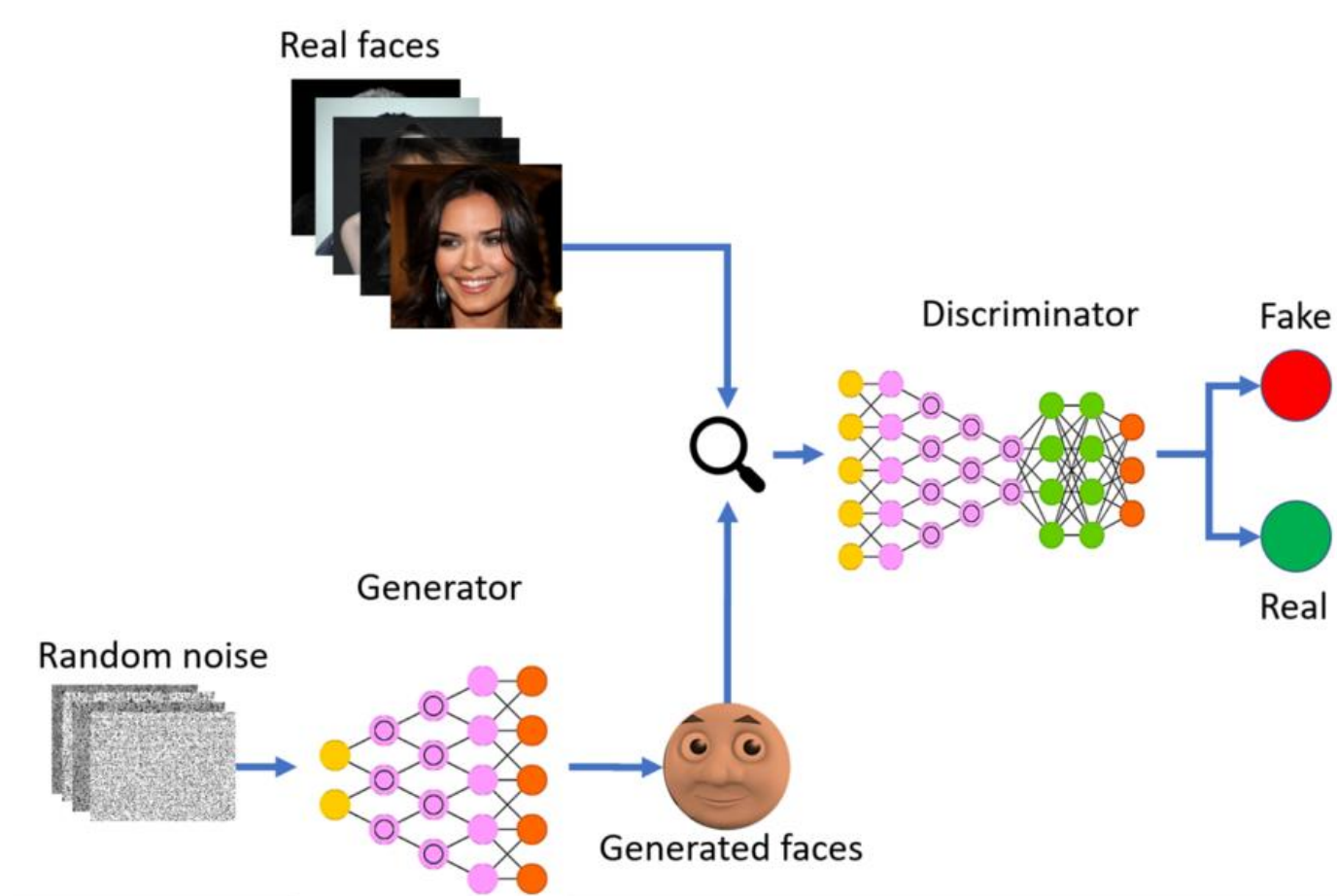
- Less widely used
- Wider range of tasks with different invariance properties

Data Augmentation GANs

- Example for emotion classification
- 5%–10% increase in the classification accuracy



GANs for faces generation



Data Augmentation for Big Data

- Can increase diversity
- Improve robustness

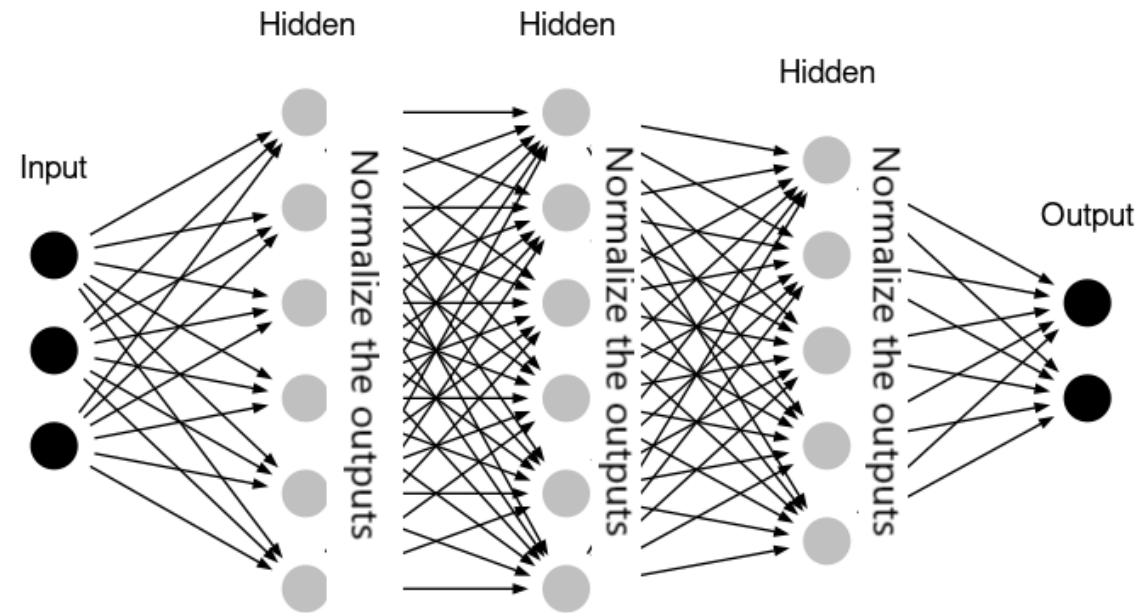
Be careful!

- Not every transformation ok!



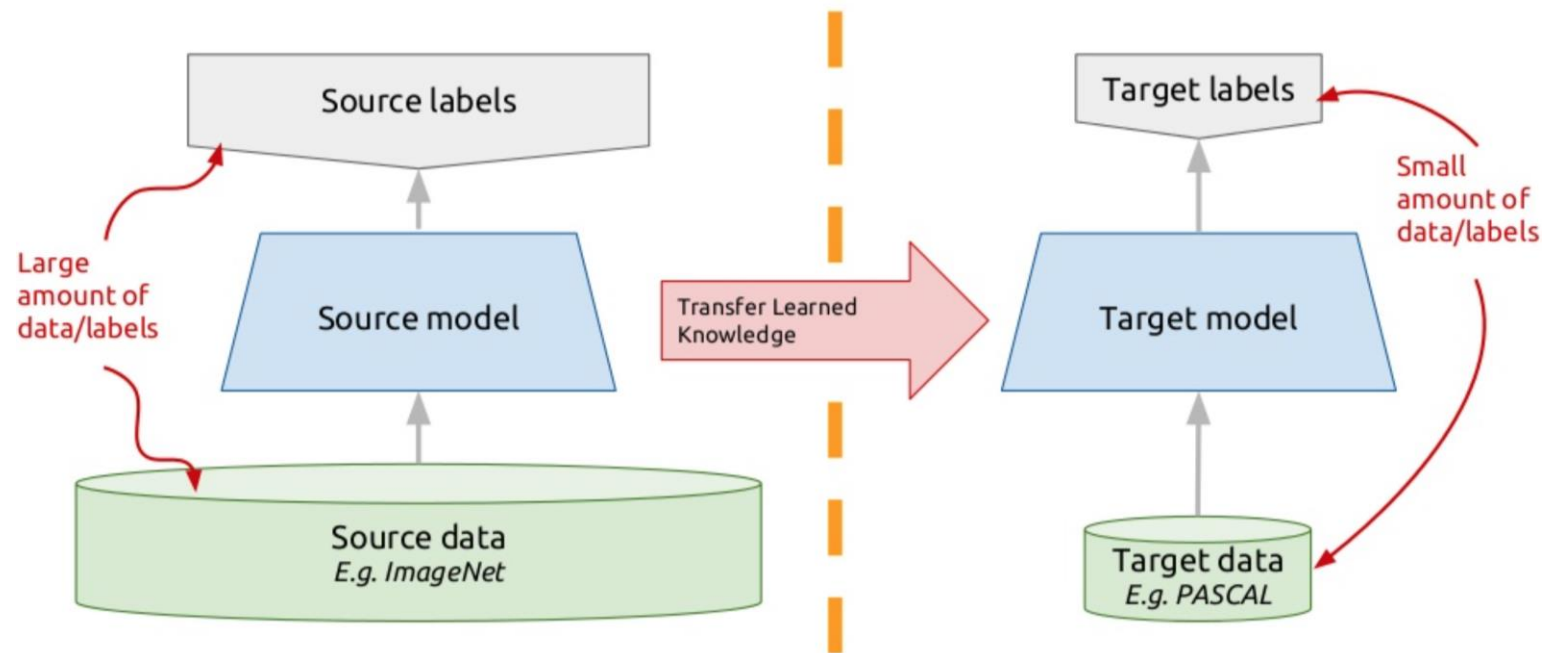
Other tools

- Batch Normalization
 - Normalize layer inputs by subtracting the batch mean and dividing by the batch standard deviation



Other tools

- Transfer Learning
 - Reuse parts of a previously trained model on a new network to solve a different but similar problem



References

Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep learning. MIT press.

Friedman, J., Hastie, T., & Tibshirani, R. (2001). The elements of statistical learning (Vol. 1, No. 10). New York: Springer series in statistics.

Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on Image Data Augmentation for Deep Learning. Journal of Big Data, 6(1), 1-48.

<https://towardsdatascience.com/understanding-the-bias-variance-tradeoff-165e6942b229>

<http://scott.fortmann-roe.com/docs/BiasVariance.html>