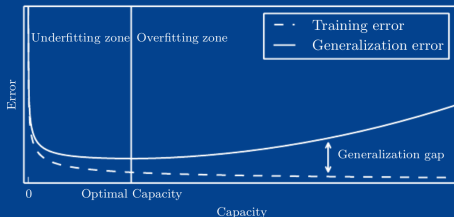# LESSON 8: Model-capacity, Under- and Overfitting, Generalization

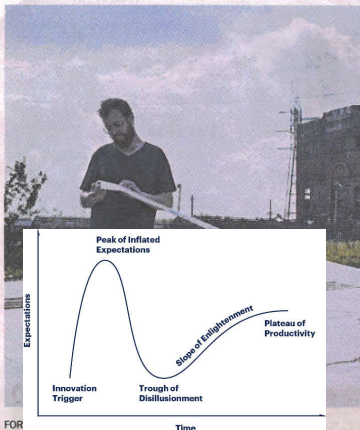## CARSTEN EIE FRIGAARD

AUTUMN 2020

# På roadtrip med en insekthjerne

**Af MIKKEL BORIS**

I 2016 slog computeren AlphaGo den 18-dobbelte verdensmester i brætspillet Go, Lee Sedol. Go er et kompliceret og abstrakt spil, som kræver intuition og kreativitet, men den kunstige intelligens vandt med en række innovative træk overlegent.

Undervejs i spillet troede kommentatorerne

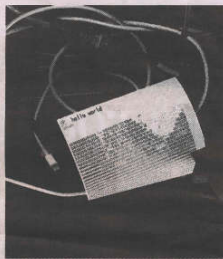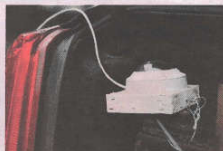Goodwin til Weekendavisen fra sin lejlighed i Los Angeles.

Inden køreturen havde han brugt måneder på at træne maskinen. Han satte den til at læse et stort korpus af moderne litteratur fra hele verden, så den kunne lære at skrive af de store forfattere.

»Det fungerer ligesom autokorrekturen på din telefon, bare klogere og trænet på et mere litterær skue. Den skriver bogstav for bogstav, så den har lært sig selv at forudsige det næste

når du dekonstruerer dem. Efter at have læst den i ét stræk og fået turen lidt på afstand har romanen fået en universalitet, så jeg kan projicere mine egne oplevelser ind i teksten,« uddyber Goodwin.

– Du har beskrevet projektet som at lære en insekthjerne at skrive. Hvad betyder det?

»Jeg forsøgte at pointere, at maskinen ikke er på niveau med den menneskelige hjerne. Et artificielt neuralt net er en algoritme, der er lavet, som vi tror, hjerner fungerer. Jeg synes





»Det er et forsøg på at skabe en ny bruger-flade for at skrive. På en måde har jeg skrevet en roman med en bil,« fortæller Ross Goodwin om sin AI-forfattede bog *1 the road*.

e Count: 8
e Count: 0
A    T
[0 0]
[0 0]
[0 0]

side

front

side

front

front

# A computer vision system to monitor the infestation level of Varroa destructor in a honeybee colony
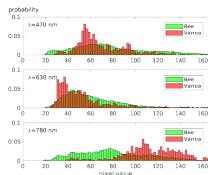


Figure 5: Histograms of bee and varroa pixel intensity values, for the spectral wavelengths 470 nm, 630 nm, and 780 nm respectively, recorded with the JAI camera. The image path is via the mirror-window-mirror, i.e. data were sampled with the setup given in figure 1. The image data for the histogram is the single bee with mite seen in figure 6.
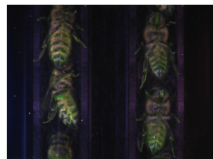


Figure 6: The actual unprocessed camera view of the bees

spectively. The CM analysis was able to rank all wavelengths combinations, using one, two, three or four district wavelengths to give a ranking list of 'best' combination also taking the JAI camera spectrum into account.

The CM value of the actual choose wavelengths combination (470-630-780 nm) gave a rank just below the CM average score. This CM analysis was conducted after picking the actual used wavelengths, so later versions of the VMU might want to investigate a CM combination with a higher rank.

A specially designed diffuser and a number of narrow spectral LED were mounted in the camera foca diffuse illuminati flections.

Figure 6 disp along the passa era, with the gre with the NIR m

### 2.3.3. Real-time processin

A color and tion of 1296×96 from the camer over two separat sary sustained b ing frames real-

These data w line post-proces first matching t rally coalescing producing a 24 the later image
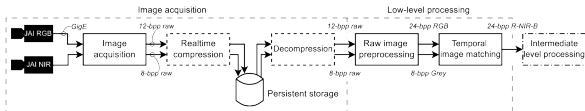
Lossless real can be applied bandwidth than



Figure 7: The low-level image processing pipeline. Raw camera images are stored on disk for later retrieval and post-processing. 12- and 8 bits-per-pixel are used as the raw JAI/Bayer packed pixel format for the RGB and IR images respectively. Lossless, real-time compression can be introduced if persistent storage bandwidth is less than the raw-stream image rate of 93 MiB/sec. The 12- and 8-bpp raw images from the network arrives out-of-order with respect to each other, hence the need for the temporal image matching.
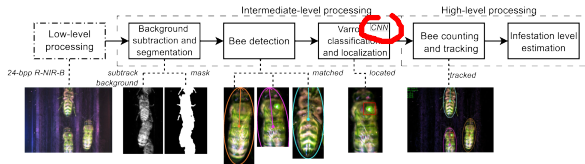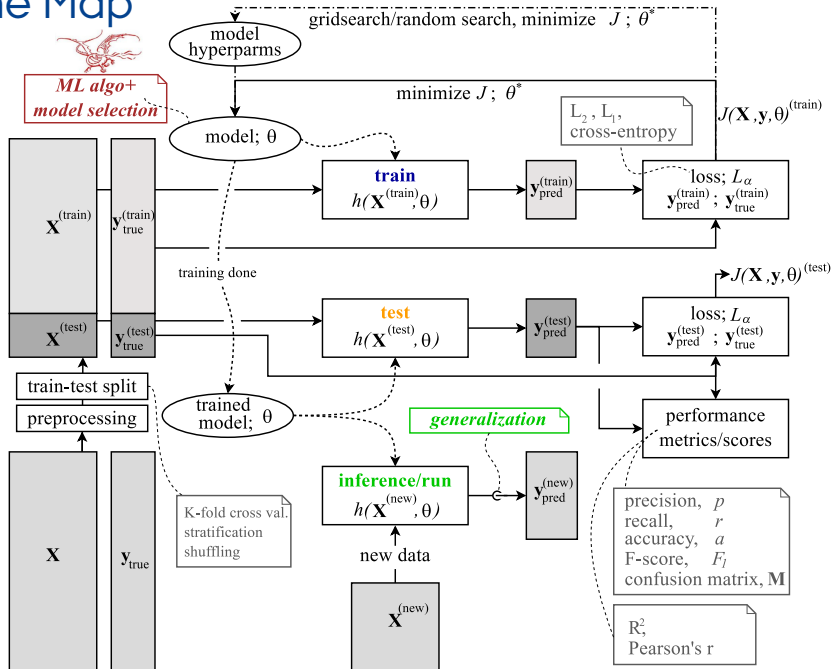


Figure 8: The processing pipeline of the intermediate- to high-level image processing algorithms to analyze and count the number of bees with Varroa destructor. A trained convolutional neural network (CNN) was used for the Varroa classification and localization stage.

# BA Project: generic tagging/labling tool

# The Map

# Preprocessing of Data
## Scaling, Standardization, Normalization...

Why the need for preprocessing?

> *Standardization of datasets is a* **common requirement for many machine learning estimators** *[..]; they might* **behave badly** *if the individual features do not more or less look like standard normally distributed data. [..]*
>
> *[https://scikit-learn.org/stable/modules/preprocessing.html]*

Standardization of a feature vector **x**, giving **x**′ mean zero, and standard deviation one

$$\mathbf{x}' = \frac{\mathbf{x} - \mu\_\mathbf{x}}{\sigma\_\mathbf{x}}$$

What kind of estimators needs standardization?
   $\rightarrow$ **Neural networks (NNs) in particular!**

Regularization and optimization:
   $\rightarrow$ WARM-UP for the comming use of NNs!

# Pipelines

Feature: GDP per capita feature in range 10K to 50K \$.

But MLP expects input in the range [0;1] or perhaps [-1;1].

```
1   # Manual scaling..
2   X_min = np.min(X)
3   X_max = np.max(X)
4   s = X_max − X_min
5
6   print(f"X_min={X_min:.0f},  X_max={X_max:.0f},  s={s:.0f}")
7
8   X_scaled = (X−X_min)/s
9   print(f"X_scaled.shape={X_scaled.shape}")
10  print(f"np.min(X_scaled)={np.min(X_scaled)}")
11  print(f"np.max(X_scaled)={np.max(X_scaled)}")
12
13  mlp.fit(X_scaled ,y.ravel())
14  y_pred_mlp = mlp.predict((M−X_min)/s)
15
16  plt.plot(m, y_pred_lin, "r")
17  #plt.plot(m, y_pred_knn, "b")
18  plt.plot(m, y_pred_mlp, "k")
19
20  print(f"mpl.score={mlp.score(X_scaled, y.ravel()):0.2f}")
```

```
Prints:
X_min=9055,
X_max=55805, s=46750
X_scaled.shape=(29, 1)
np.min(X_scaled)=0.0
np.max(X_scaled)=1.0
mpl.score=0.70
```

# Pipelines

Housing Data and MLPs: introducing a `MinMaxScaler`

```python
1   # Now, do the same but via a pipeline..
2   from sklearn.preprocessing import MinMaxScaler
3
4   scaler = MinMaxScaler()
5   scaler.fit(X)
6   X_scaled2 = scaler.transform(X)
7   M_scaled = scaler.transform(M)
8
9   mlp.fit(X_scaled2, y.ravel())
10  y_pred_mlp = mlp.predict(M_scaled)
11
12  print(f"mpl.score={mlp.score(X_scaled, y.ravel()):0.2f}")
13
14  # PRINTS: mpl.score=0.71
```

# Pipelines

Housing Data and MLPs: putting everything in a Full Pipeline

```
1    # Or even better, in a full pipeline..
2    from sklearn.pipeline import Pipeline
3
4    pipe = Pipeline(
5      [
6        ('scaler', MinMaxScaler()),
7        ('mlp', mlp)
8      ])
9    pipe.fit(X, y.ravel())
10
11   print(f"pipe.score(..) = {pipe.score(X, y.ravel()):0.2f}"
12
13   # PRINTS: mpl.score=0.68
```

# Pipelines

Brief intro to Scikit-learn pipelines..

Python code from `capacity_under_overfitting.ipynb`

```python
1   from sklearn.pipeline import Pipeline
2   from sklearn.preprocessing import PolynomialFeatures
3   from sklearn.linear_model import LinearRegression
4   from sklearn.model_selection import cross_val_score
5
6   ..
7
8   polynomial_features = PolynomialFeatures(degree=degrees[i], ..
9   linear_regression = LinearRegression()
10
11  pipeline = Pipeline([
12          ("polynomial_features", polynomial_features),
13          ("linear_regression", linear_regression)
14      ])
15  pipeline.fit(X[:, np.newaxis], y)
16
17  scores = cross_val_score(pipeline, X[:, np.newaxis], y, scoring=
        "neg_mean_squared_error", cv=10)
18  ..
19  score_mean = -scores.mean()
```

## Classification metrics

See the Classification metrics section of the user guide for further details.

| | |
|---|---|
| `metrics.accuracy_score` (y_true, y_pred[, …]) | Accuracy classification score. |
| `metrics.auc` (x, y[, reorder]) | Compute Area Under the Curve (AUC) using the trapezoidal rule |
| `metrics.average_precision_score` (y_true, y_score) | Compute average precision (AP) from prediction scores |
| `metrics.cohen_kappa_score` (y1, y2[, labels, …]) | Cohen's kappa: a statistic that measures inter-annotator agreement. |
| `metrics.confusion_matrix` (y_true, y_pred[, …]) | Compute confusion matrix to evaluate the accuracy of a classification |
| `metrics.f1_score` (y_true, y_pred[, labels, …]) | Compute the F1 score, also known as balanced F-score or F-measure |
| `metrics.log_loss` (y_true, y_pred[, eps, …]) | Log loss, aka logistic loss or cross-entropy loss. |
| `metrics.precision_score` (y_true, y_pred[, …]) | Compute the precision |
| `metrics.recall_score` (y_true, y_pred[, …]) | Compute the recall |
| `metrics.roc_auc_score` (y_true, y_score[, …]) | Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores. |
| `metrics.roc_curve` (y_true, y_score[, …]) | Compute Receiver operating characteristic (ROC) |
| `metrics.zero_one_loss` (y_true, y_pred[, …]) | Zero-one classification loss. |

**Notes on Keras MLPs**

Typical Keras MLP Supervised Classifier setup.

▶ loss function
```
loss='categorical_
       crossentropy'
```
▶ metrics collected via history
```
metrics=[
    'categorical_accuracy',
    'mean_squared_error',
    'mean_absolute_error'])
```
▶ input lay.: categorical encoding.
▶ output lay.: softmax function.

And notice that Keras do *not* provide metrics like
```
precision, recall, F1
```
but instead
```
categorical_accuracy, binary_accuracy
```

## Regression metrics

See the Regression metrics section of the user guide for further details.
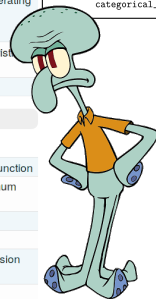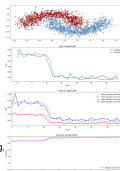
| | |
|---|---|
| `metrics.explained_variance_score` (y_true, y_pred) | Explained variance regression score function |
| `metrics.max_error` (y_true, y_pred) | max_error metric calculates the maximum residual error. |
| `metrics.mean_absolute_error` (y_true, y_pred) | Mean absolute error regression loss |
| `metrics.mean_squared_error` (y_true, y_pred[, …]) | Mean squared error regression loss |
| `metrics.mean_squared_log_error` (y_true, y_pred) | Mean squared logarithmic error regression loss |
| `metrics.median_absolute_error` (y_true, y_pred) | Median absolute error regression loss |
| `metrics.r2_score` (y_true, y_pred[, …]) | R^2 (coefficient of determination) regression |

# Model capacity

Dummy and Paradox classifier:
*capacity* fixed $\sim 0$, cannot generalize at all!

Linear regression for a polynomial model:
*capacity* $\sim$ degree of the polynomial, $x^n$

Neural Network model:
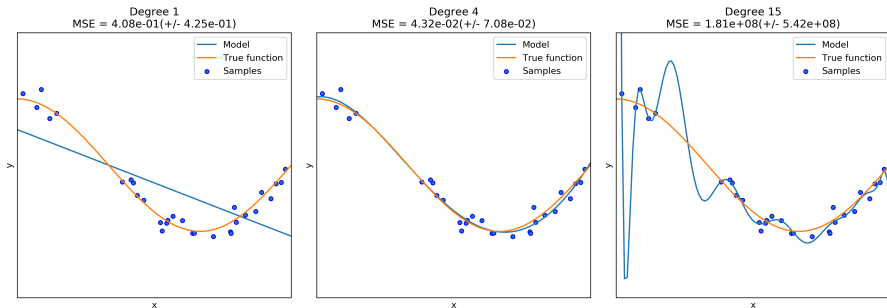*capacity* $\propto$ number of neurons/layers

$\Rightarrow$ **Capacity** can be hard to express as a quantity for some models, but you need to choose..

$\Longrightarrow$ how to choose the **optimal** *capacity*?

# Under- and overfitting

Polynomial linear reg. fit for underlying model: `cos(x)`



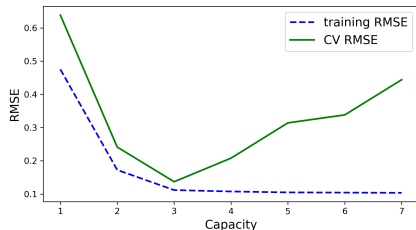▶ underfitting: capacity of model too low,
▶ overfitting: capacity to high.

⟹ how to choose the **optimal** capacity?
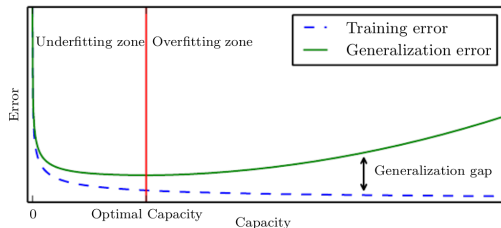
# Generalization Error

Exercise: `generalization_error.ipynb`

RMSE-capacity plot for lin. reg. with polynomial features

(capacity $\sim$ degree of poly)

(Figure 5.3 from [DL])



Inspecting the plots from the exercise (`.ipynb`) and [DL], extracting the concepts:

- ► training/generalization error,
- ► generalization gab,
- ► underfit/overfit zone,
- ► optimal capacity (best-model, early stop),
- ► (and the two axes: x/capacity, y/error.)

# Generalization Error

Exercise: `generalization_error.ipynb`

NOTE: three methods/plots:
- i) via **learning curves** as in [HOML],
- ii) via an **error-capacity** plot as in [GITHOML] and [DL],
- iii) via an **error-epoch** plot as in [GITHOML].