# Adversarial examples, face verification
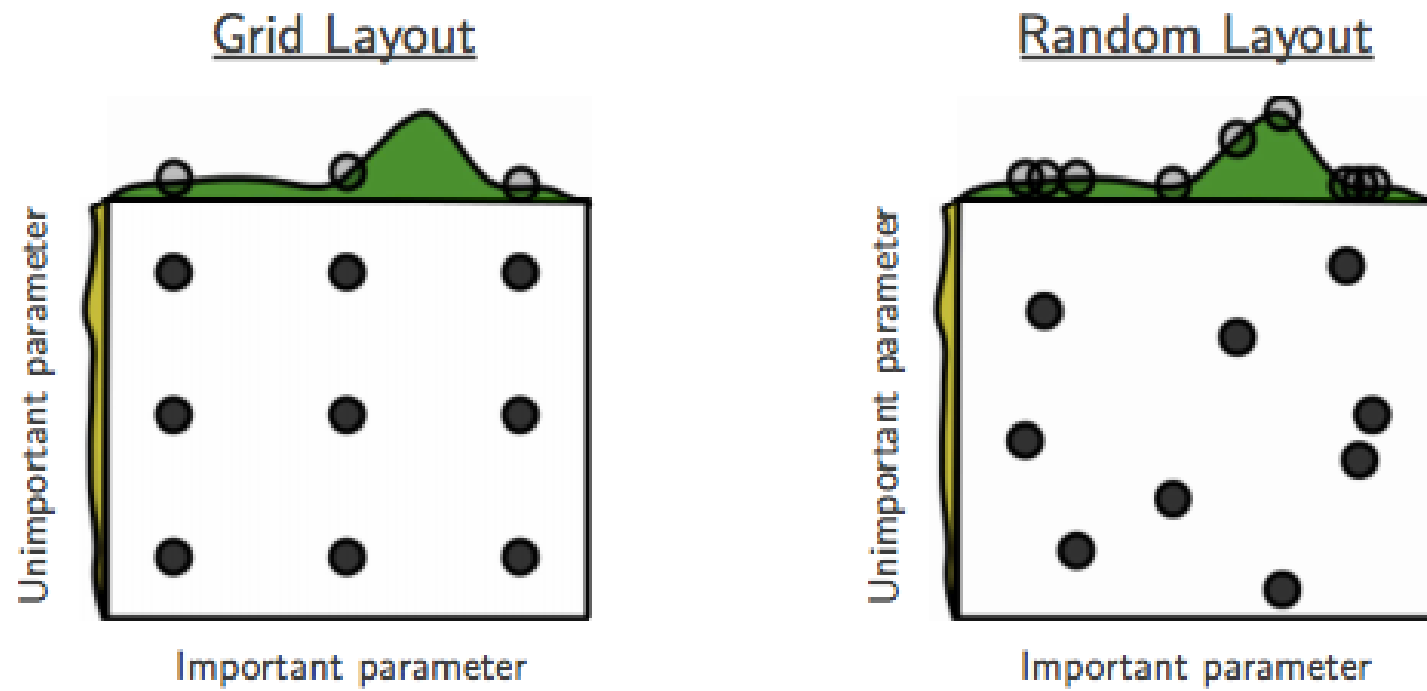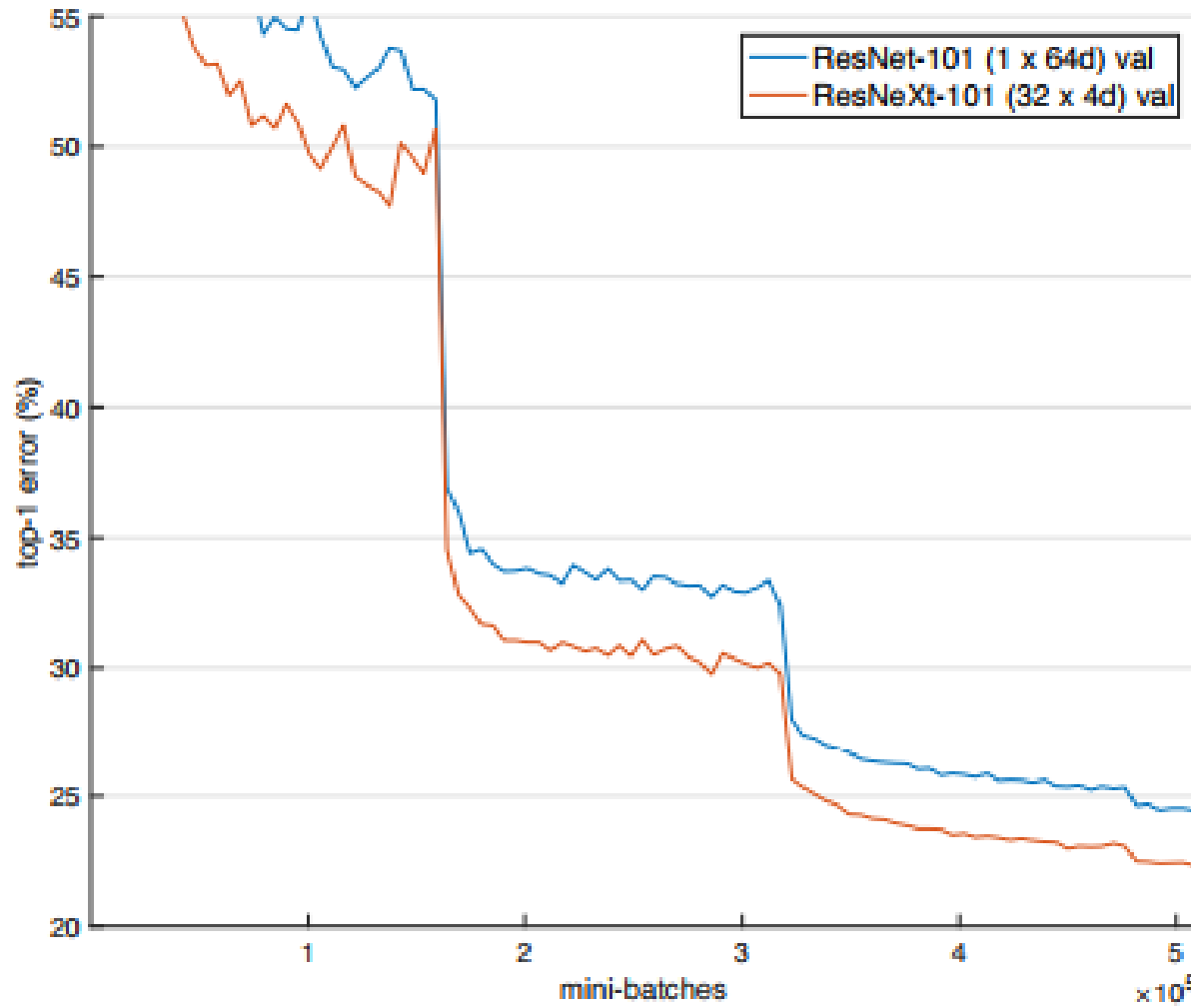
Bálint Ármin Pataki

**Hyperparameter**: the one are tuned by hand
**Parameter**: the ones are tuned by the optimization algorithm

Diagrammatic representation of Grid Search by Bergstra & Bengio
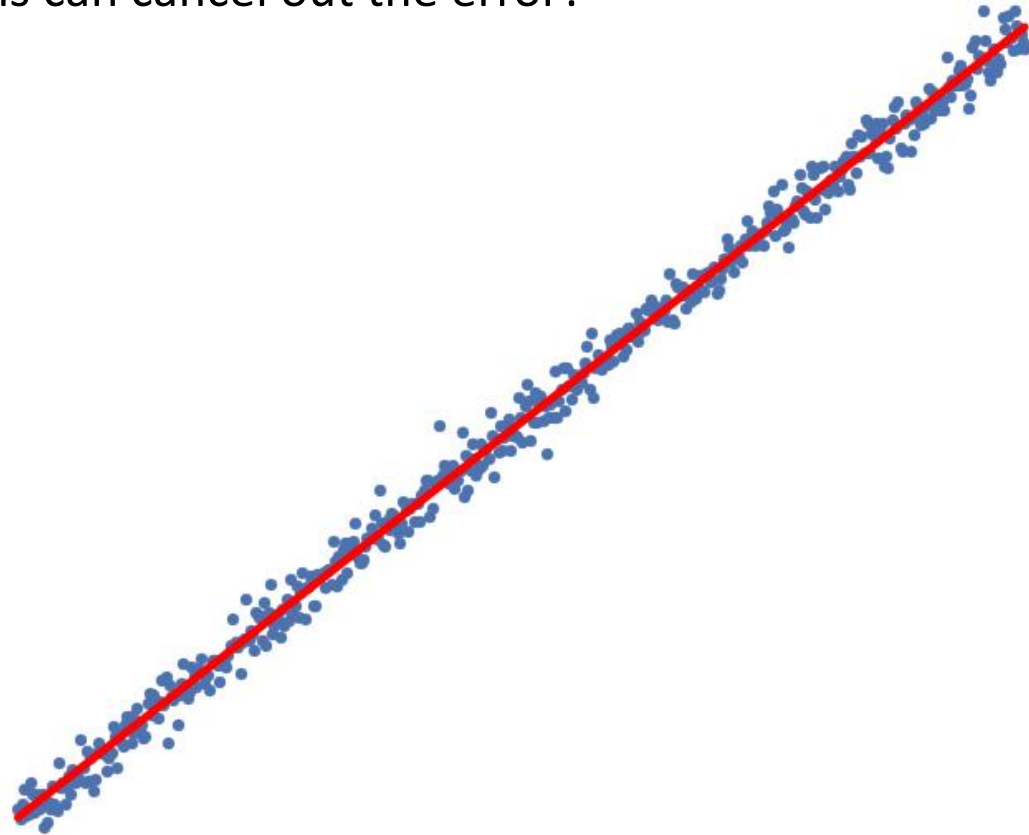


Grid Layout — Random Layout

## Pre-defined schedule vs baby-sitting



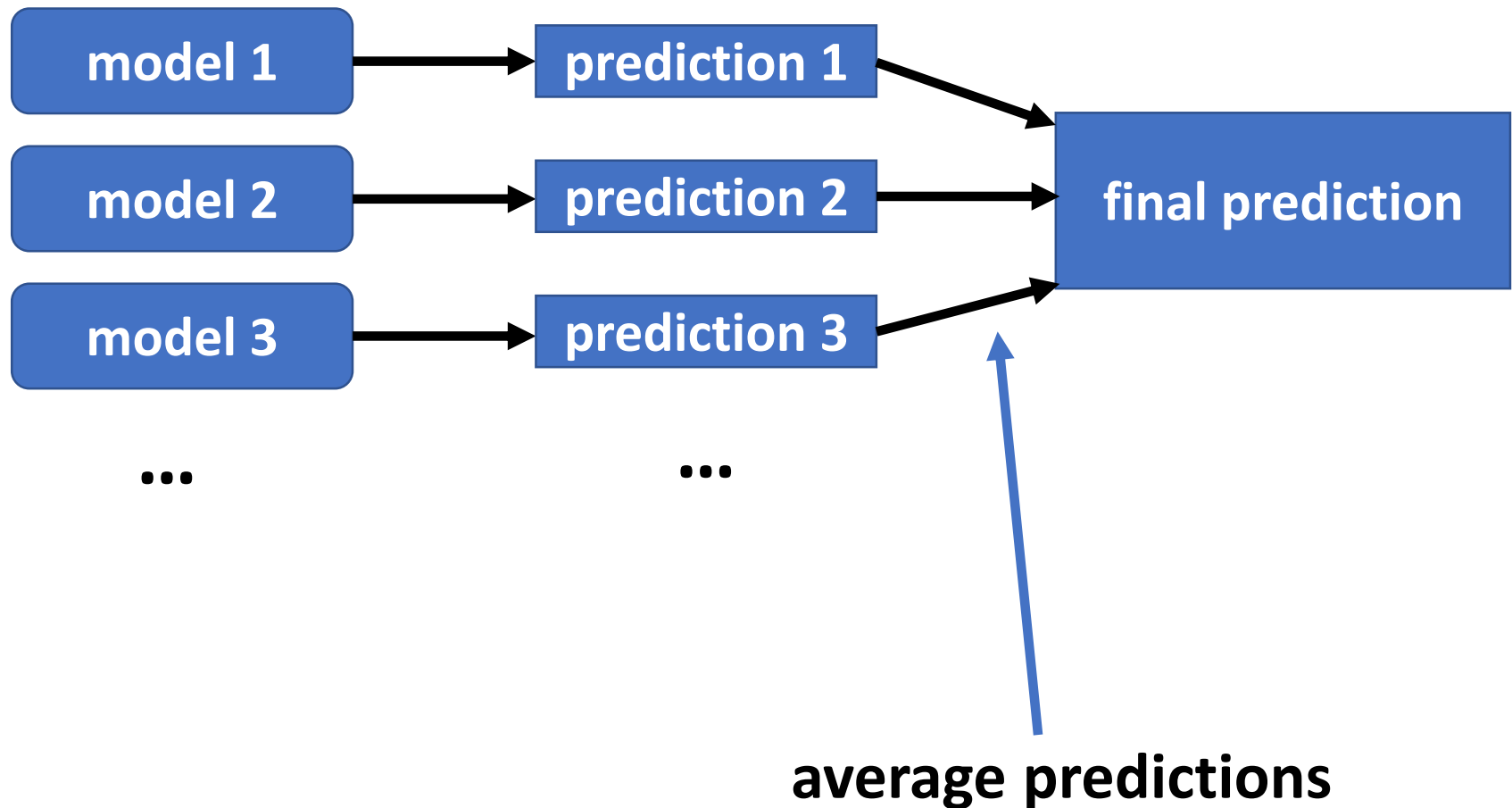[Xie: Aggregated Residual Transformations for Deep Neural Networks  arXiv:1611.05431v2]

What if the error is random, but model dependent?
→training different models can cancel out the error?
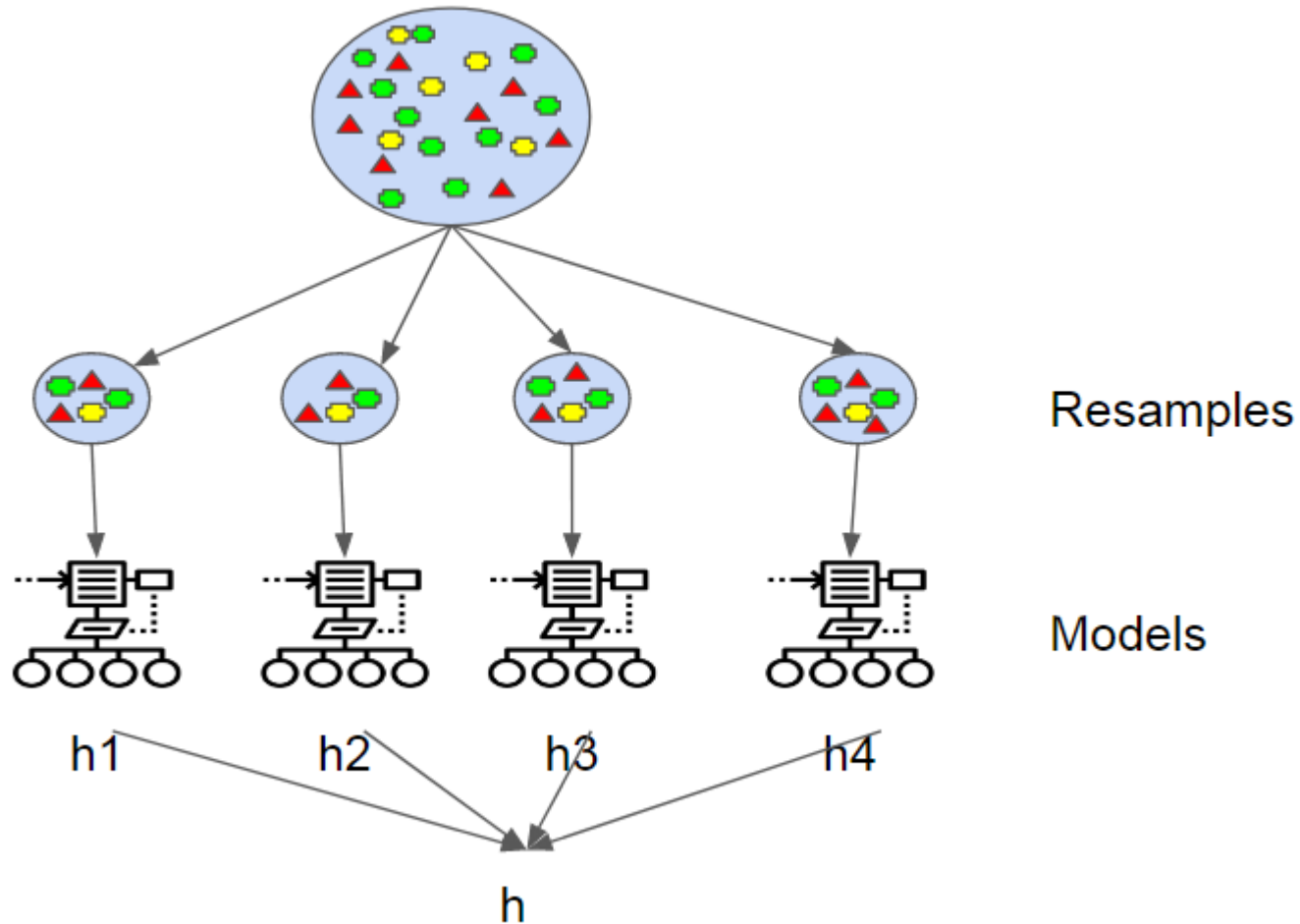
At Kaggle always an ensemble wins.
Often it has no practical relevance, but it can increase the score with epsilon.

**average predictions**

Special case: Bagging (Bootstap AGGregatING)
- same models trained on subset of train data

Resamples

Models

h1    h2    h3    h4

h

https://medium.com/@SeattleDataGuy/how-to-develop-a-robust-algorithm-c38e08f32201

Model 1 fits the original data.
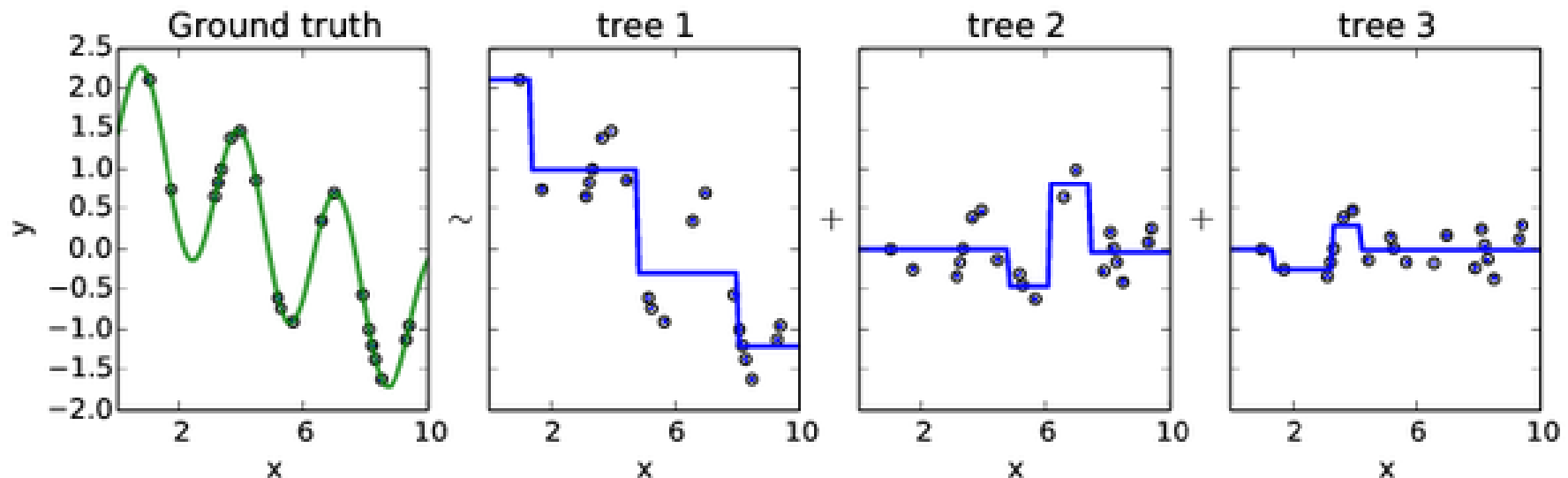
Model 2 fits original data – model 1 prediction = the residual

Model 3 fits original data – model 1 prediction – model 2 prediction

…



https://www.quora.com/How-would-you-explain-gradient-boosting-machine-learning-technique-in-no-more-than-300-words-to-non-science-major-college-students

Model 1 is string for smaller X, model 2 is for larger X.
Would be great to combine them!

# Stacking

```
model 1  →  prediction 1 ⎤
model 2  →  prediction 2 ⎥→  model  →  final prediction
model 3  →  prediction 3 ⎦
```

...          ...



model1
model2

# Proper validation strategy



- **diff. view of same event**
- **diff. measurement of same person**
- **time series**

- **...**

# DEMO notebook

Deep learning and machine learning in science

# Other PPT

http://www.robots.ox.ac.uk/~vgg/publications/2015/Parkhi15/presentation.pptx

# Importance of the correct metric

https://www.walesonline.co.uk/news/wales-news/facial-recognition-wrongly-identified-2000-14619145

Facial recognition software wrongly identified more than 2,000 people as potential criminals as police patrolled the Champions League final in Cardiff.

The technology provided hundreds of "false positives" wrongly marking out innocent people as possible troublemakers when an estimated 170,000 people descended on the city for the showpiece match between Real Madrid and Juventus.

A South Wales Police spokesman admitted "no facial recognition system is 100% accurate under all conditions" but added that in the months since it was first deployed "no-one has been arrested where a 'false positive alert' has occurred and no members of the public have complained".

Data published by the force showed police covering the Champions League final at the Principality Stadium on June 3 last year were alerted to 2,470 potential matches with custody pictures by the facial recognition programme.

But of these 92% – a total if 2,297 – were incorrect, with just 173 providing 'true positive alerts'.

MNIST

 - one vs all classifier (zero or not zero)

 - metric: accuracy

 - 90% accuracy. Is it good?

# Instead of accuracy

Positive = predicted class

True = prediction is correct →

| True Positive (TP) | True Negative (TN) |
|---|---|
| False Positive (FP) | False Negative (FN) |

Is this digit 0 (actually it is)?
Prediction: yes
→True positive

Is this digit 0 (actually it is not)?
Prediction: no
→True negative

Is this digit 0 (actually it is not)?
Prediction: yes
→False positive

Is this digit 0 (actually it is)?
Prediction: no
→False negative

You predict a probability of being positive.
Then a threshold is applied (eq 50%) and is the probability is above, then prediction is positive

For different thesholds there is different prediction.

For different tasks you want different goals:
 - identification of criminals: avoid False Negative
 - unlock your phone with face recog: avoid False Positive

http://www.navan.name/roc/



Receiver operating characteristic example

https://stats.stackexchange.com/questions/132777/what-does-auc-stand-for-and-what-is-it

Problem dependent. Objectives:
- easy to understand/interpret
- significantly better model should have significantly better score
- cover your exact need (as possible)

| | | True condition | | | |
|---|---|---|---|---|---|
| | Total population | Condition positive | Condition negative | Prevalence = $\frac{\Sigma \text{ Condition positive}}{\Sigma \text{ Total population}}$ | Accuracy (ACC) = $\frac{\Sigma \text{ True positive} + \Sigma \text{ True negative}}{\Sigma \text{ Total population}}$ |
| Predicted condition | Predicted condition positive | **True positive**, Power | **False positive**, Type I error | Positive predictive value (PPV), Precision = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Predicted condition positive}}$ | False discovery rate (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Predicted condition positive}}$ |
| | Predicted condition negative | **False negative**, Type II error | **True negative** | False omission rate (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Predicted condition negative}}$ | Negative predictive value (NPV) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Predicted condition negative}}$ |
| | | True positive rate (TPR), Recall, Sensitivity, probability of detection = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$ | False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$ | Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$ | Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$ $F_1$ score = $\frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$ |
| | | False negative rate (FNR), Miss rate = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$ | True negative rate (TNR), Specificity (SPC) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$ | Negative likelihood ratio (LR−) = $\frac{FNR}{TNR}$ | |

https://en.wikipedia.org/wiki/Precision_and_recall

**sensitivity**, **recall**, **hit rate**, or **true positive rate (TPR)**

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN}$$

**specificity** or **true negative rate (TNR)**

$$TNR = \frac{TN}{N} = \frac{TN}{TN + FP}$$

**precision** or **positive predictive value (PPV)**

$$PPV = \frac{TP}{TP + FP}$$

**negative predictive value (NPV)**

$$NPV = \frac{TN}{TN + FN}$$

**miss rate** or **false negative rate (FNR)**

$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - TPR$$

**fall-out** or **false positive rate (FPR)**

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR$$

**false discovery rate (FDR)**

$$FDR = \frac{FP}{FP + TP} = 1 - PPV$$

**false omission rate (FOR)**

$$FOR = \frac{FN}{FN + TN} = 1 - NPV$$

**accuracy (ACC)**

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$

**F1 score**

is the harmonic mean of precision and sensitivity

$$F_1 = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$

**Matthews correlation coefficient (MCC)**

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

**Informedness** or **Bookmaker Informedness (BM)**

$$BM = TPR + TNR - 1$$

**Markedness (MK)**

$$MK = PPV + NPV - 1$$

*Sources: Fawcett (2006), Powers (2011), and Ting (2011)* [4] [1] [5]

https://en.wikipedia.org/wiki/Precision_and_recall

# Importance of the correct metric

https://www.walesonline.co.uk/news/wales-news/facial-recognition-wrongly-identified-2000-14619145

Facial recognition software wrongly identified more than 2,000 people as potential criminals as police patrolled the Champions League final in Cardiff.

The technology provided hundreds of "false positives" wrongly marking out innocent people as possible troublemakers when an estimated 170,000 people descended on the city for the showpiece match between Real Madrid and Juventus.

A South Wales Police spokesman admitted "no facial recognition system is 100% accurate under all conditions" but added that in the months since it was first deployed "no-one has been arrested where a 'false positive alert' has occurred and no members of the public have complained".

Data published by the force showed police covering the Champions League final at the Principality Stadium on June 3 last year were alerted to 2,470 potential matches with custody pictures by the facial recognition programme.

But of these 92% – a total if 2,297 – were incorrect, with just 173 providing 'true positive alerts'.

https://blog.openai.com/adversarial-example-research/

https://blog.openai.com/robust-adversarial-inputs/

http://www.navan.name/roc/

http://arogozhnikov.github.io/2016/07/05/gradient_boosting_playground.html

http://www.robots.ox.ac.uk/~vgg/publications/2015/Parkhi15/presentation.pptx