

Why a highly accurate model could be not the best^[1]?



Mohamed Maher

Products Group_(SWDC), Giza Systems
<https://www.menti.com/alxf2i73h12j> (4786 8330)

Background: Probabilistic Classification



Cats



Dogs

Background: Probabilistic Classification



ML Model is trained to separate between dogs and cats by creating a **decision boundary**



Cats



Dogs

99.5% Accuracy

Background: Probabilistic Classification



Input Image



ML Model

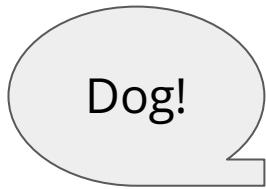


Cats



Dogs

Background: Probabilistic Classification



ML Model



Cats



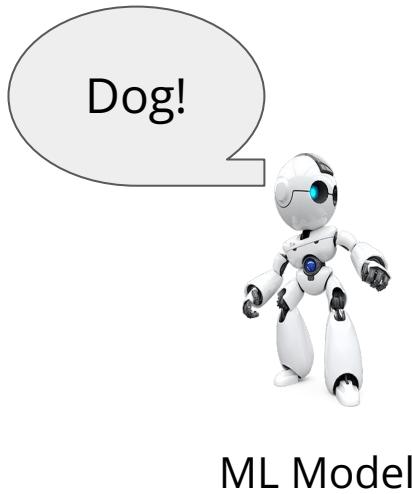
Dogs



Background: Probabilistic Classification



Background: Probabilistic Classification

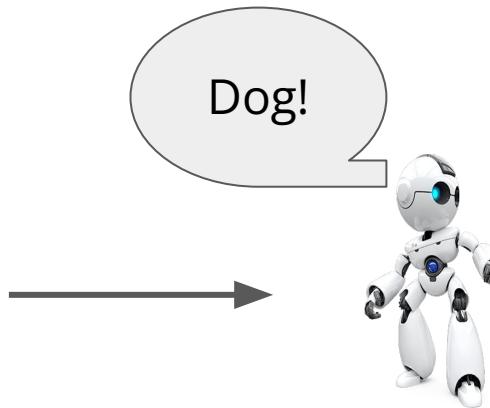
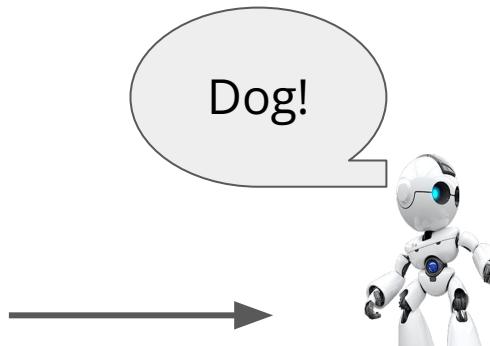


Cats



Dogs

Background: Probabilistic Classification



Background: Probabilistic Classification



Input Image

ML Model



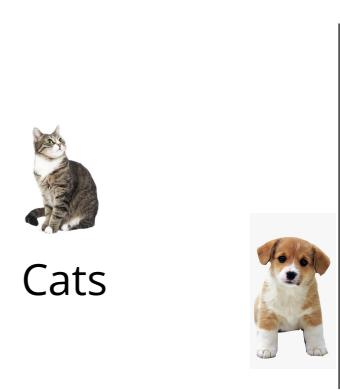
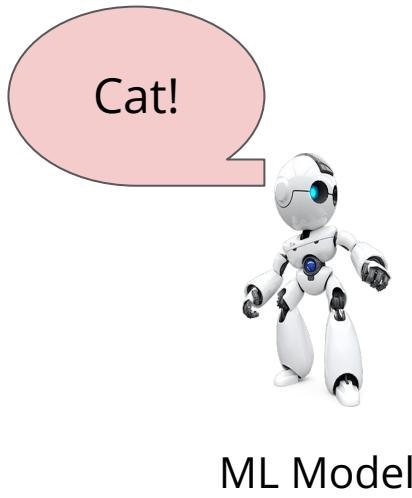
Cats



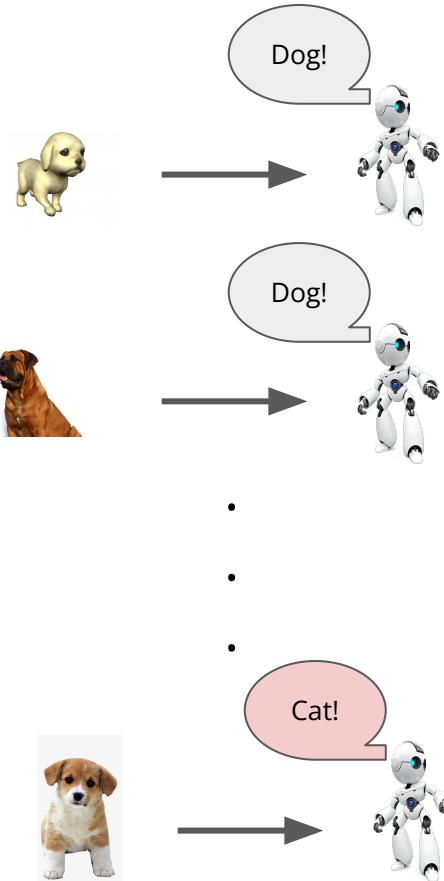
Dogs



Background: Probabilistic Classification

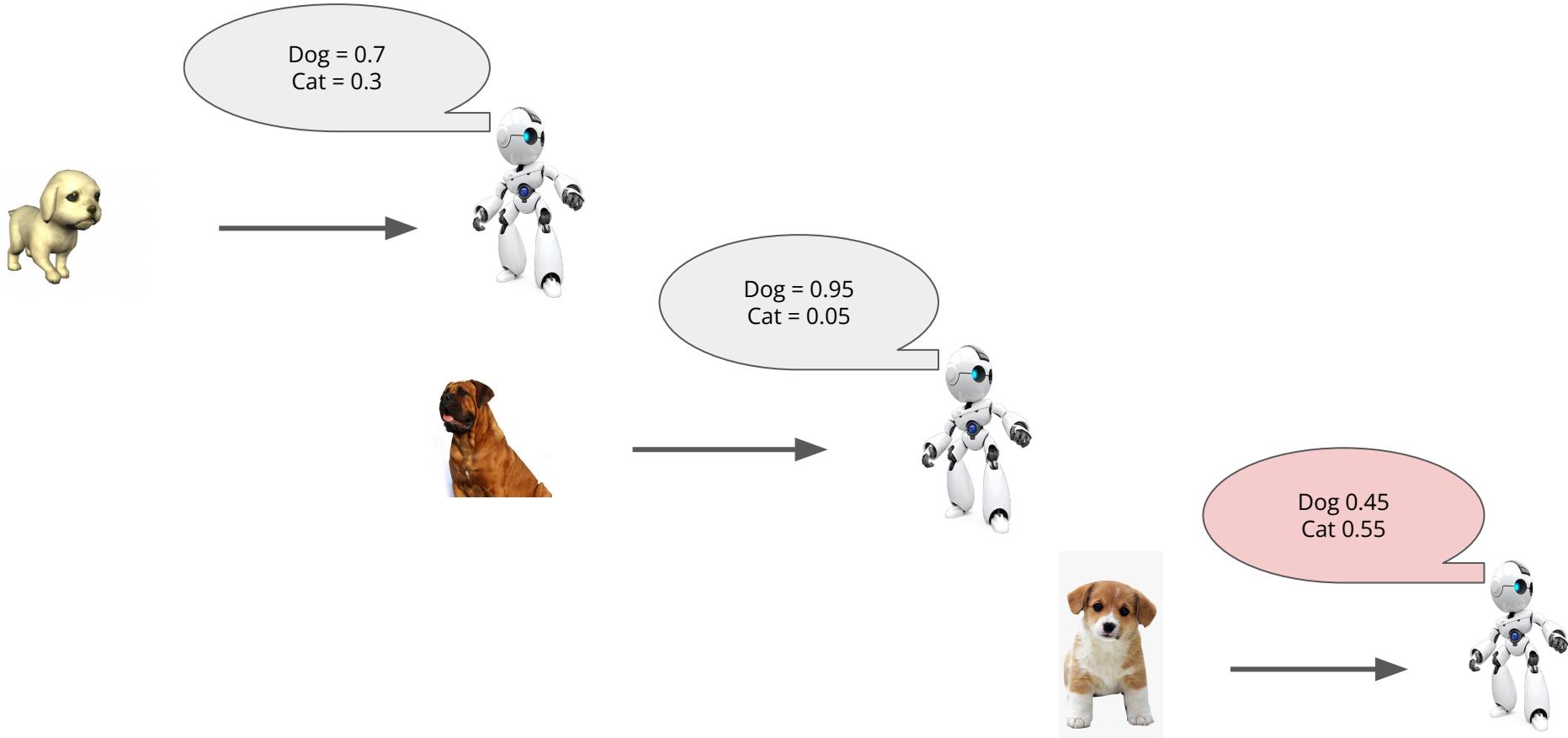


Background: Probabilistic Classification



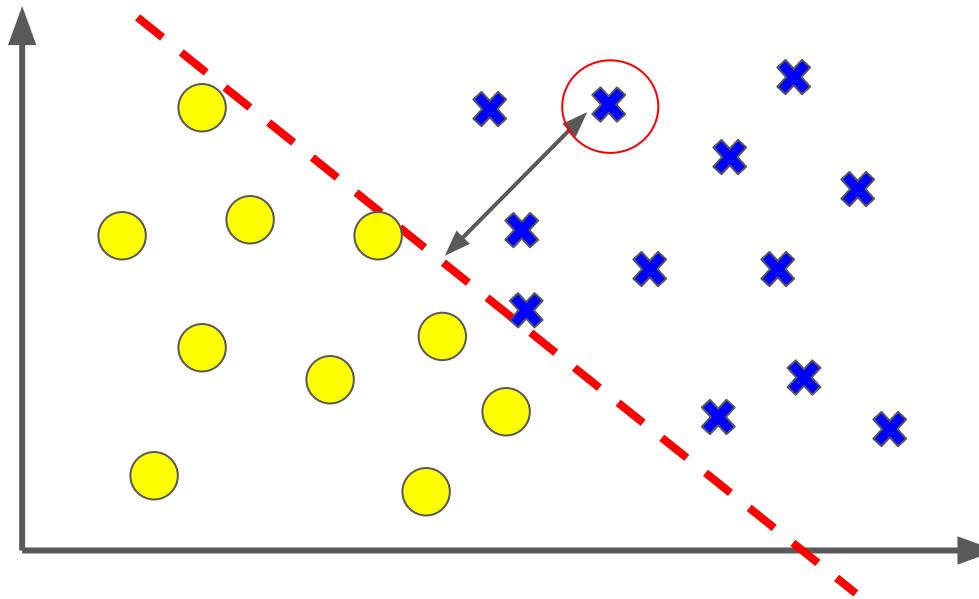
- Still 99.5% Accurate.
- In sensitive domains, I can not tolerate many wrong decisions especially when data distribution changed.
- Probabilistic Classification is needed.

Background: Probabilistic Classification



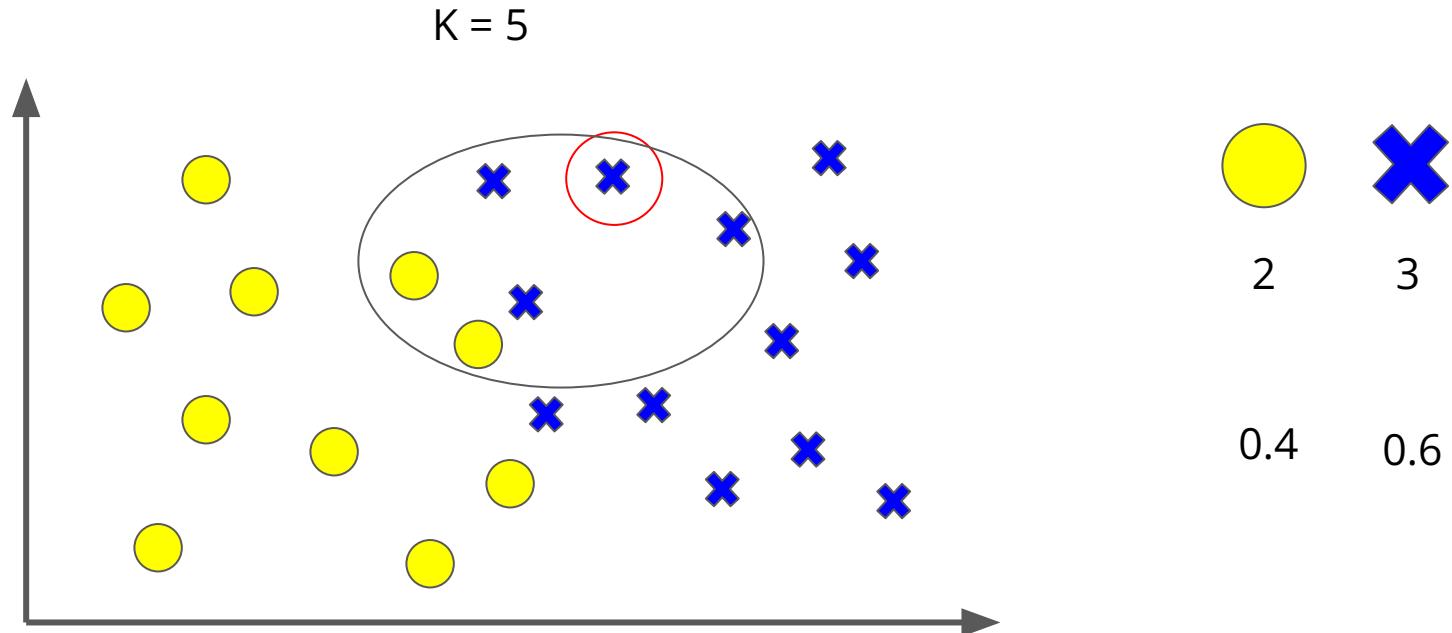
Background

- Can the model predict probabilities instead of labels only?
 - Linear models: distance to the decision boundary



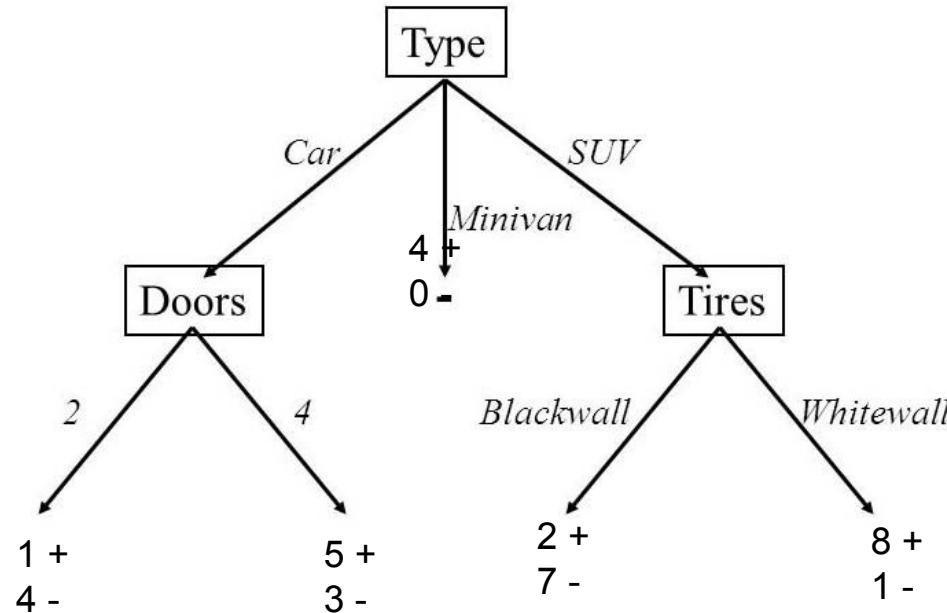
Background

- Can the model predict probabilities instead of labels only?
 - KNN: number of neighbors from one class / total number of neighbors



Background

- Can the model predict probabilities instead of labels only?
 - Decision Tree



SUV with Blackwall tires has probability $2/(2+7) = 0.22$ to be a taxi!

Motivation



AI-based diagnostic tools for COVID-19

- More than 36 published studies are proposing different machine learning models to detect the viral infection based on tomography scans ¹.
- Can not be used in practice because of 1) Interpretability 2) Calibration

[1] Wynants Laure, et al. "Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal." *bmj* 369 (2020).

The image is a collage of academic publications and journal covers. At the top left, a 'Radiology' journal cover features a study by Lin Li et al. titled 'Artificial intelligence distinguishes COVID-19 from other acute respiratory infections using chest CT scans'. The study was published in a letter on 19 May 2020. Below it, a 'Nature Medicine' cover shows a study by Xueyan Mei et al. titled 'Artificial intelligence diagnosis of patient COVID-19 pneumonia using chest CT scans'. In the center, a 'Deep Learn.' study by Xuehai He et al. is shown, titled 'Sample-Efficient Deep Learning for COVID-19 Diagnosis Based on CT Scans'. This study is published in the 'IEEE Transactions on Medical Imaging'. To the right, a 'Covid-19' study by Pengtao Xie et al. is mentioned, which includes 'IoU^{3,f}'.

Motivation



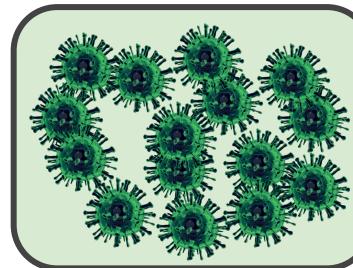
Let's use AI to diagnose COVID-19 faster

Example: Two models are used to diagnose **100** infected similar samples with COVID-19.

Model 1

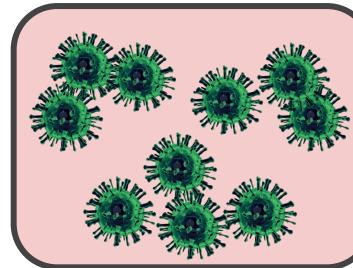


True Positives(60)



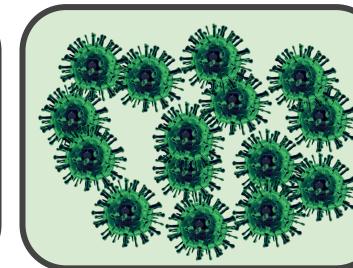
Confidence = **0.6**,

False Negatives(40)

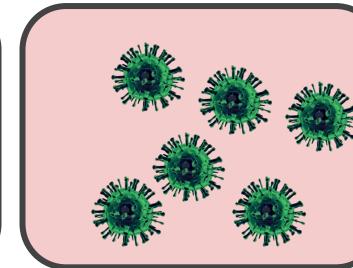


Accuracy = **60%**

True Positives(75)



False Negatives(25)



Confidence = **0.8**,

Model 2



Accuracy = **75%**

Motivation



Let's use AI to diagnose COVID-19 faster

Example: Two models are used to diagnose **another 100** infected similar samples

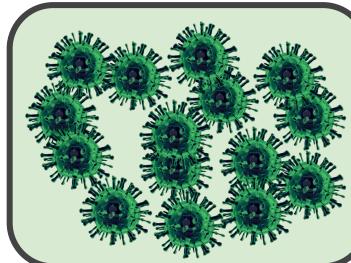
Model 1



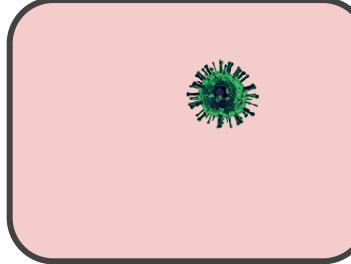
Accuracy = **99%**

Confidence = **0.99**

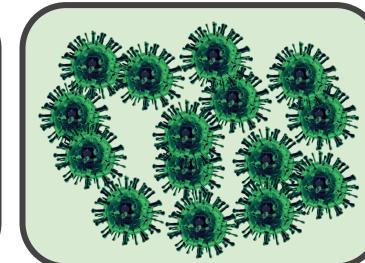
True Positives(99)



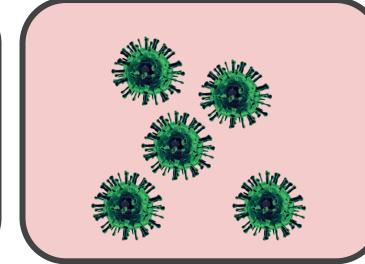
False Negatives(1)



True Positives(95)



False Negatives(5)



Model 2



Accuracy = **95%**

Confidence = **0.99**

Motivation



Which model to trust?

Model **1** has lower accuracy but if it predicts with high confidence, I would **trust** it.



Overall Accuracy
79.5 %

Model 1

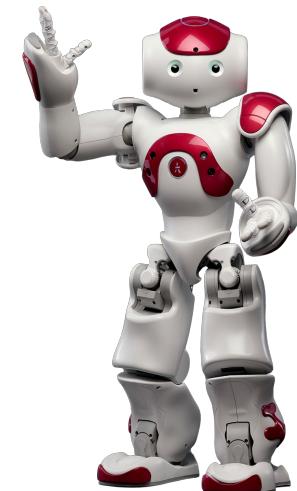
- Less accuracy
- Calibrated

If model **1** is confident of its predictions, it is almost correct.
If it is not, we can verify the results using slower standard diagnostic tests.

Model 2

- Higher accuracy
- Overconfident

If model **2** is confident of its predictions, few false negatives can lead to ***disastrous spread of the virus***



Overall Accuracy
85 %

Motivation: Calibrated Classifier

Perfect Calibrated Classifier

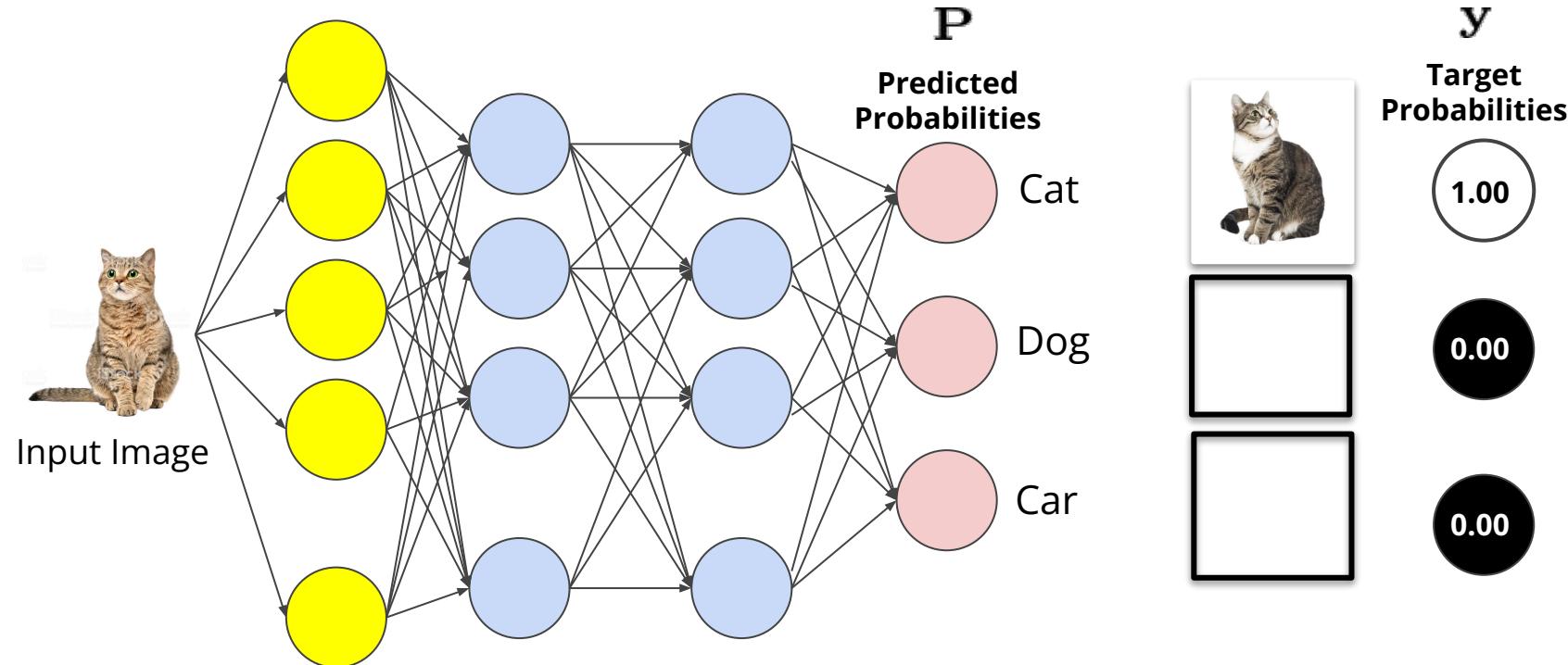
A model is *perfectly calibrated* if, for any probability value p , a prediction of a class with confidence p is correct $100*p\%$ of the time.

Model 1 is the
WINNER!



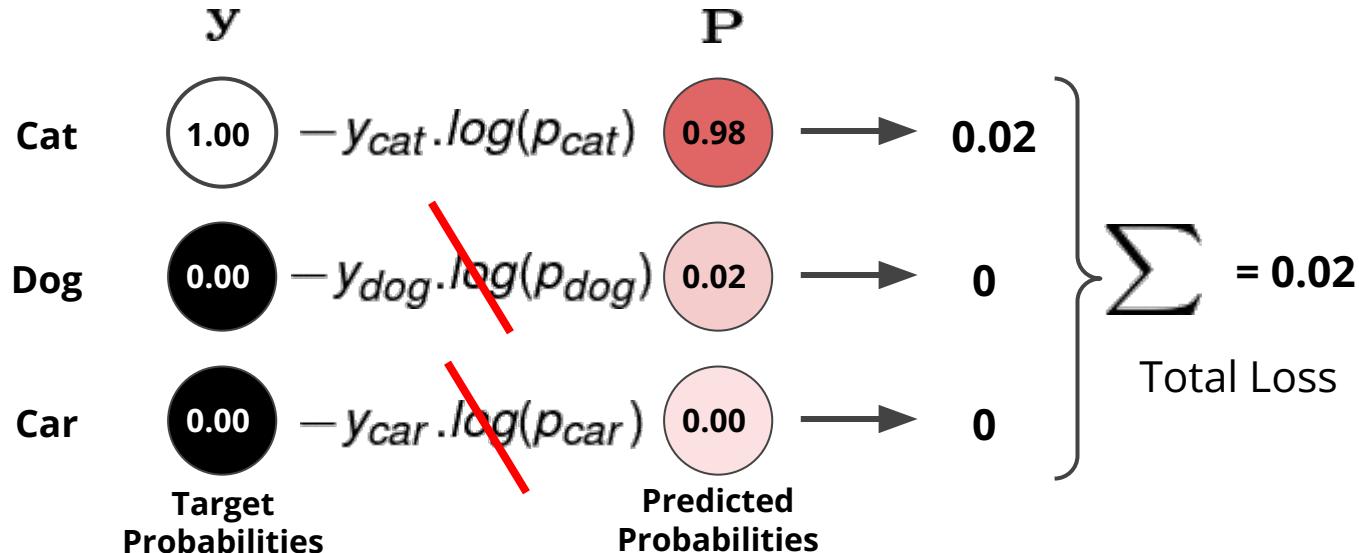
Motivation: E.g. OverConfidence in Deep Learning

Teach to predict a cat with full confidence.



Motivation: E.g. OverConfidence in Deep Learning

- Overconfident predictions



Motivation: E.g. OverConfidence in Deep Learning



CIFAR10 → Model Training

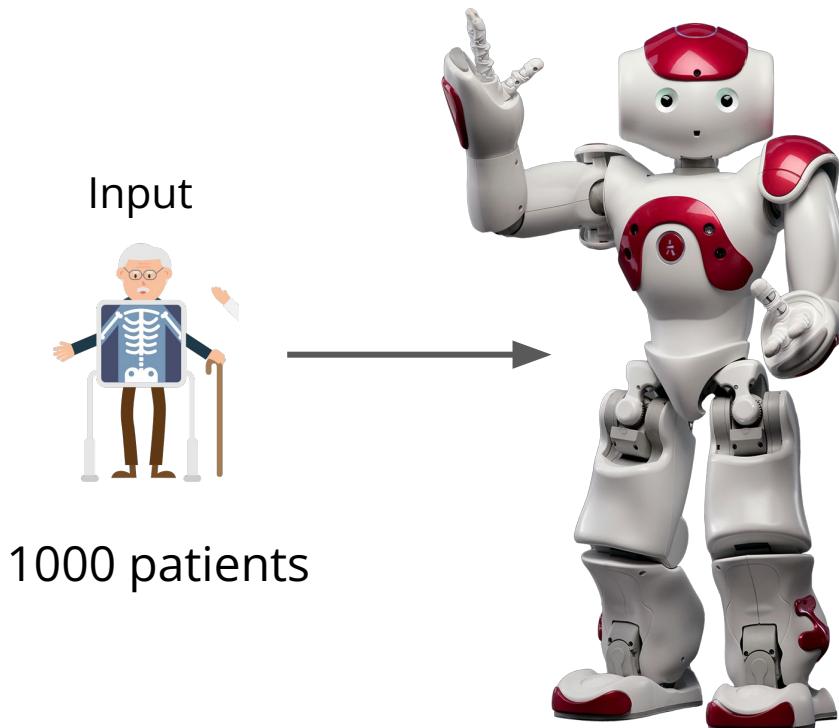


SVHN → Model Evaluation

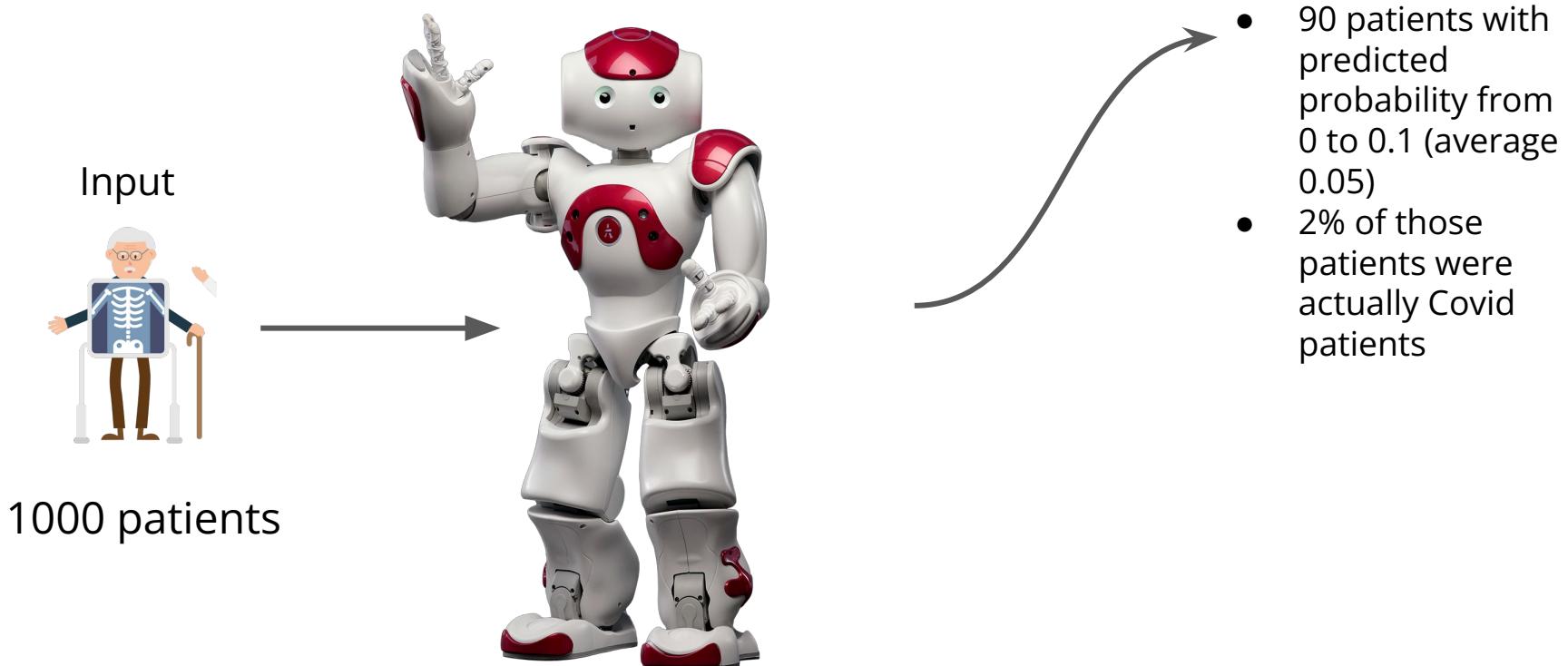
Neural Networks yield overconfident predictions on unrelated tasks!

Hein, Matthias, Maksym Andriushchenko, and Julian Bitterwolf. "Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2019.

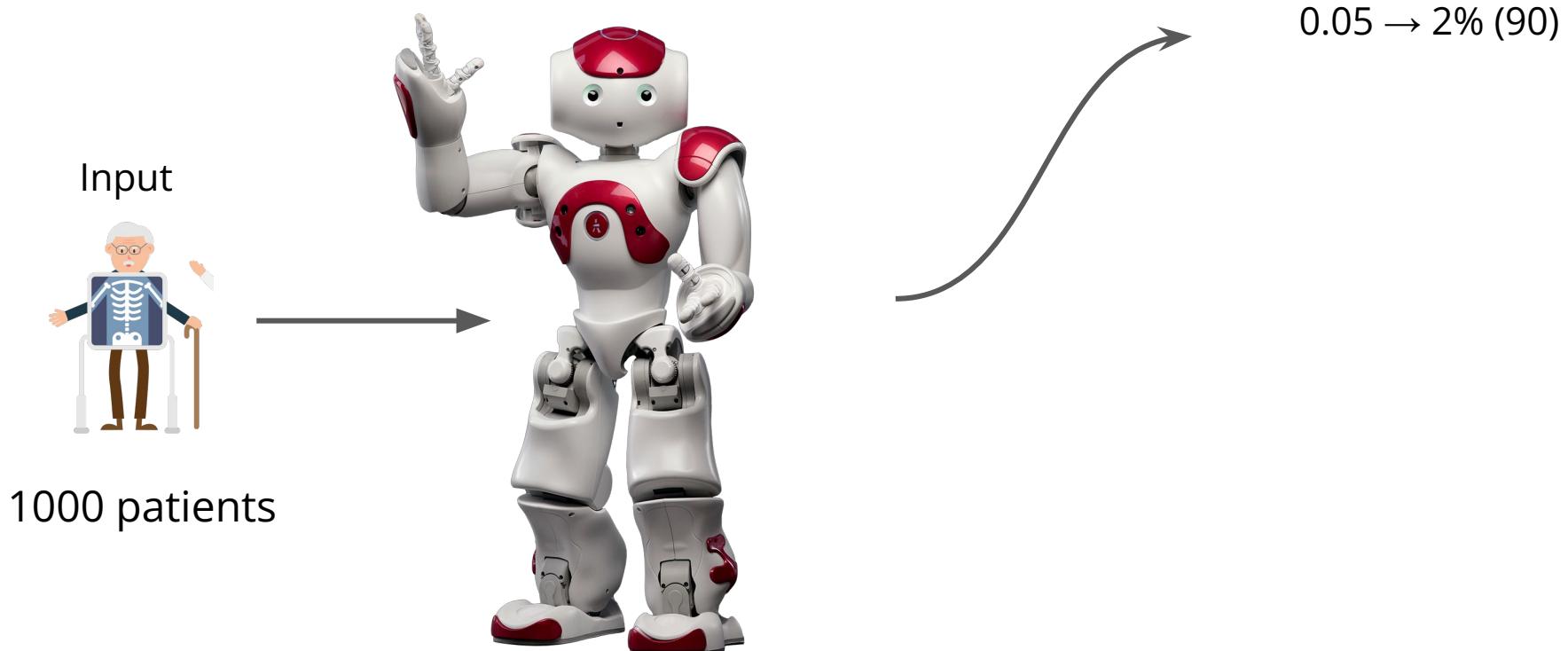
How to measure Calibration? Reliability Diagram



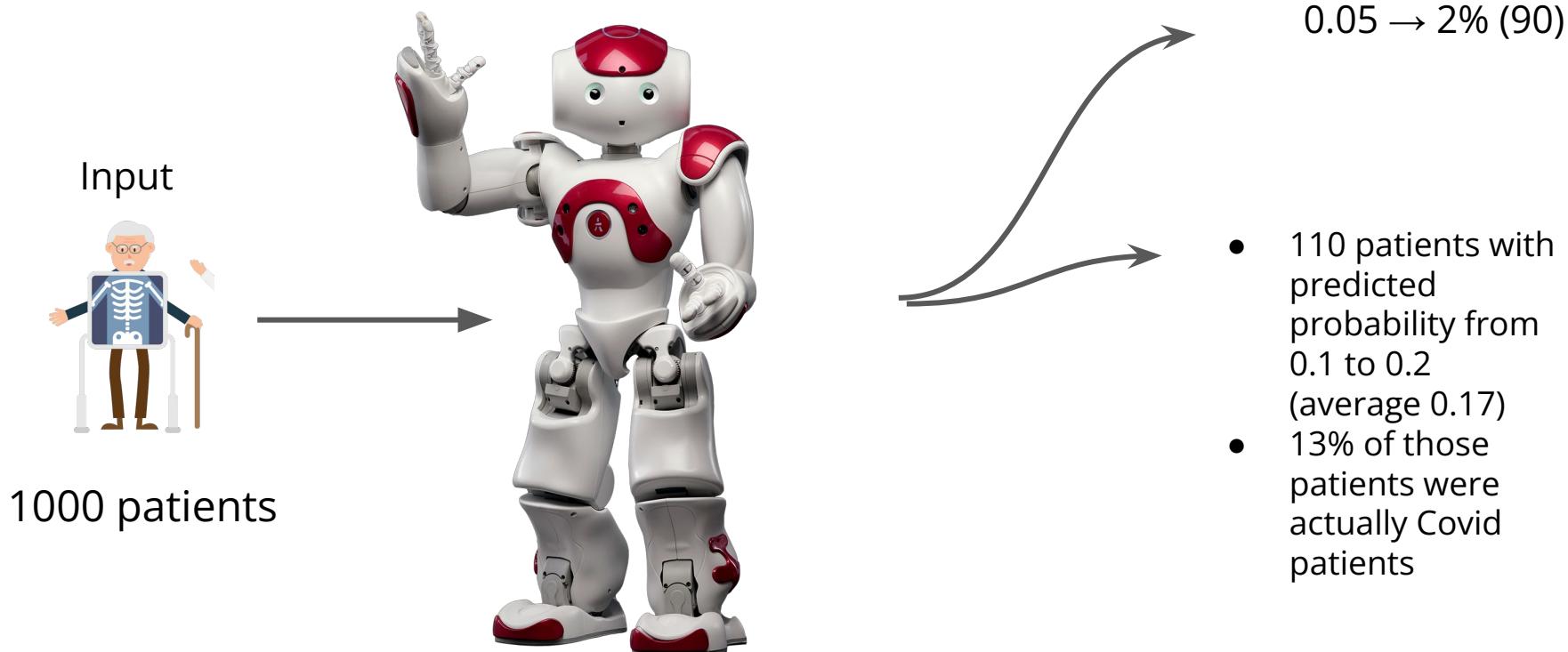
How to measure Calibration? Reliability Diagram



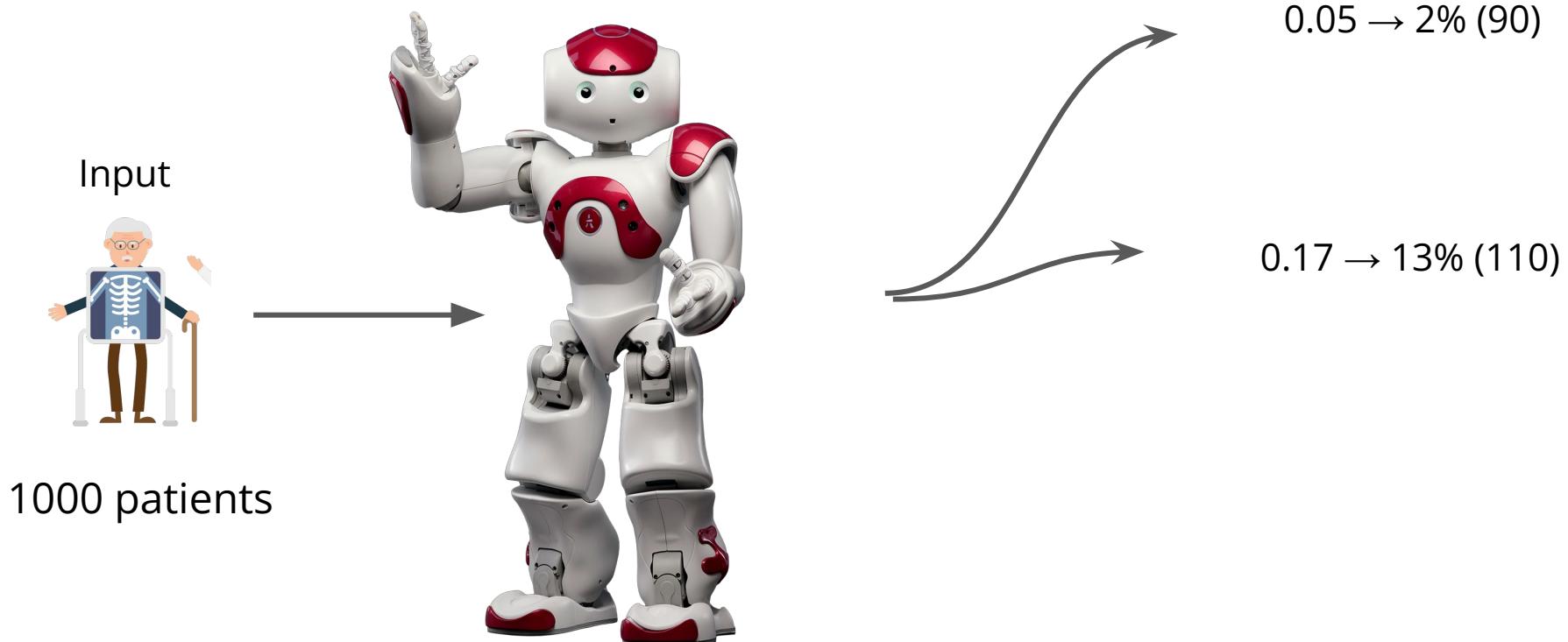
How to measure Calibration? Reliability Diagram



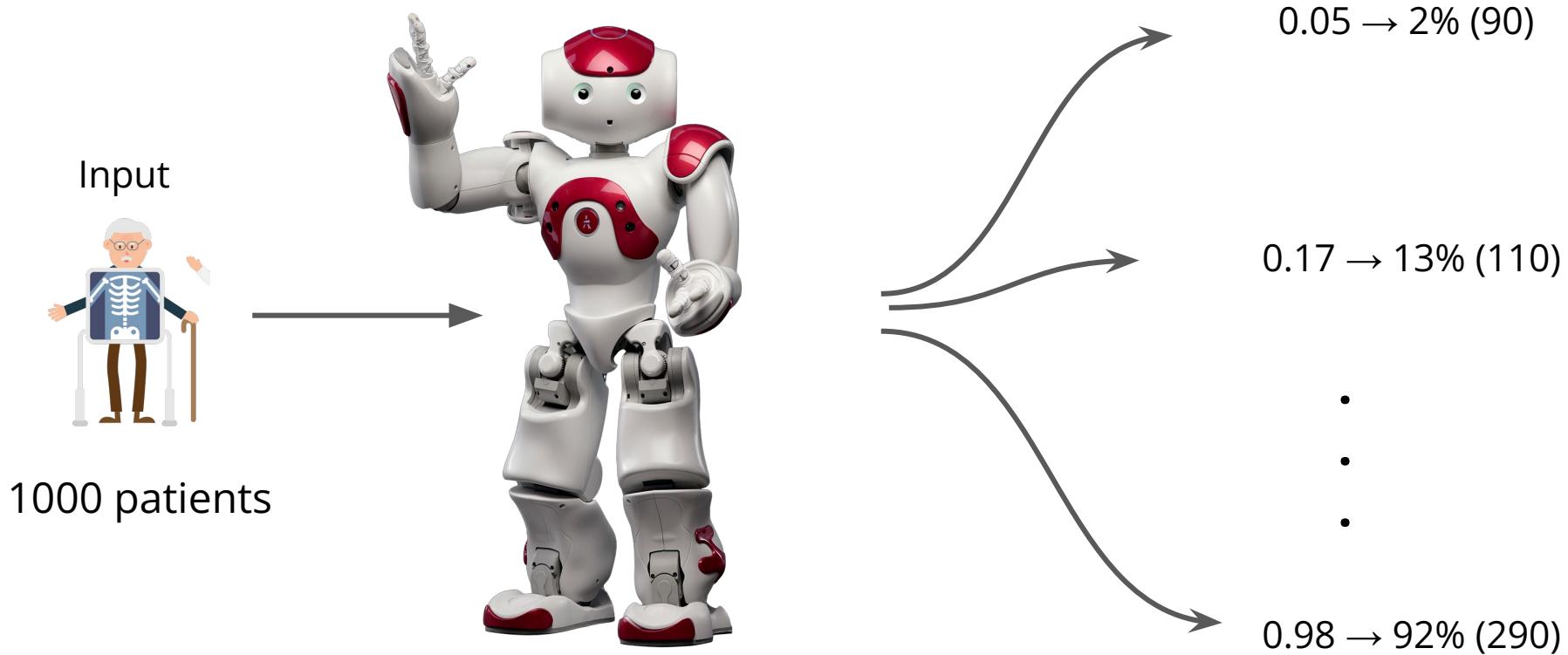
How to measure Calibration? Reliability Diagram



How to measure Calibration? Reliability Diagram

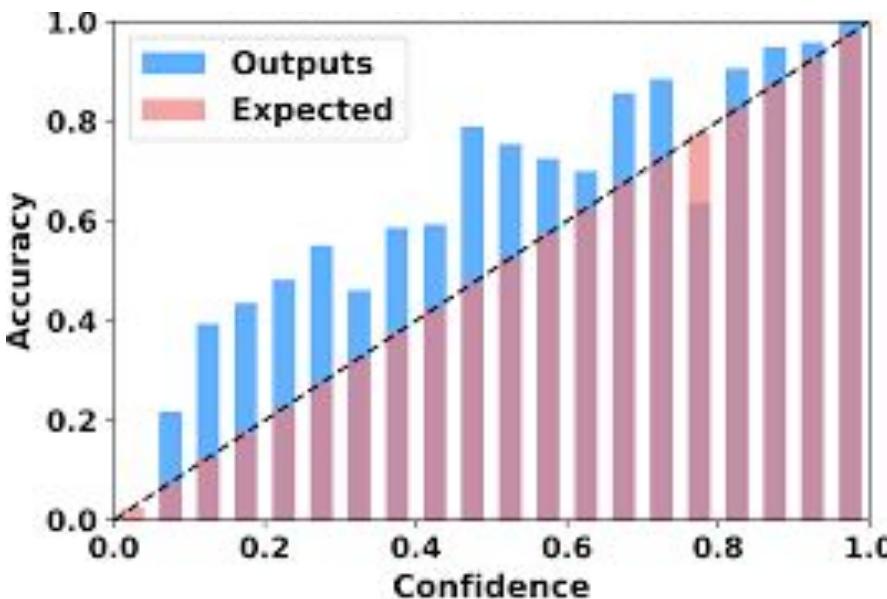


How to measure Calibration? Reliability Diagram

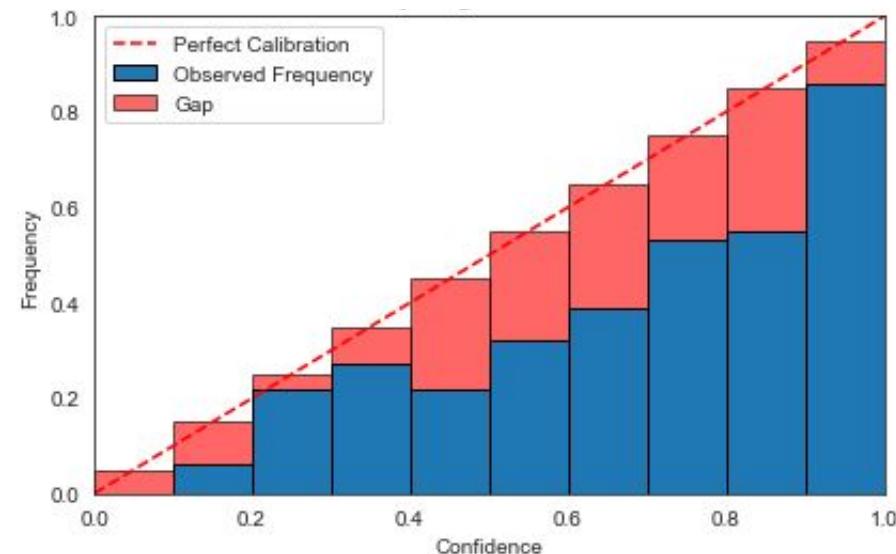


How to measure Calibration? Reliability Diagram

Underconfident Uncalibrated Model



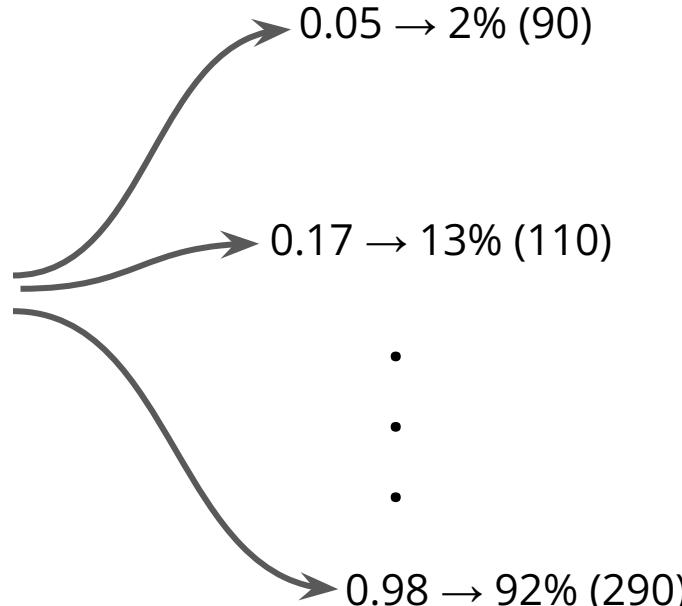
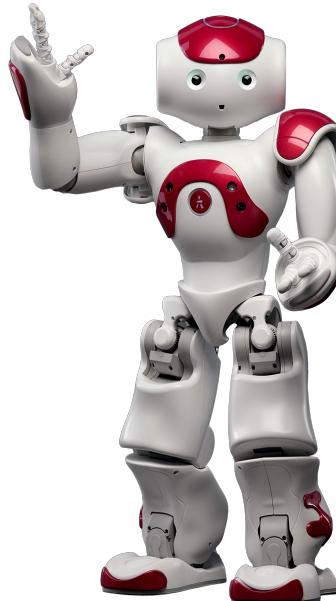
Overconfident Uncalibrated Model



How to measure Calibration? ECE

- Expected Calibration Error

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} \left| \text{acc}(B_m) - \text{conf}(B_m) \right|$$



$$90/1000 * |0.05-0.02|$$

+

$$110/1000 * |0.17-0.13|$$

+

•

•

•

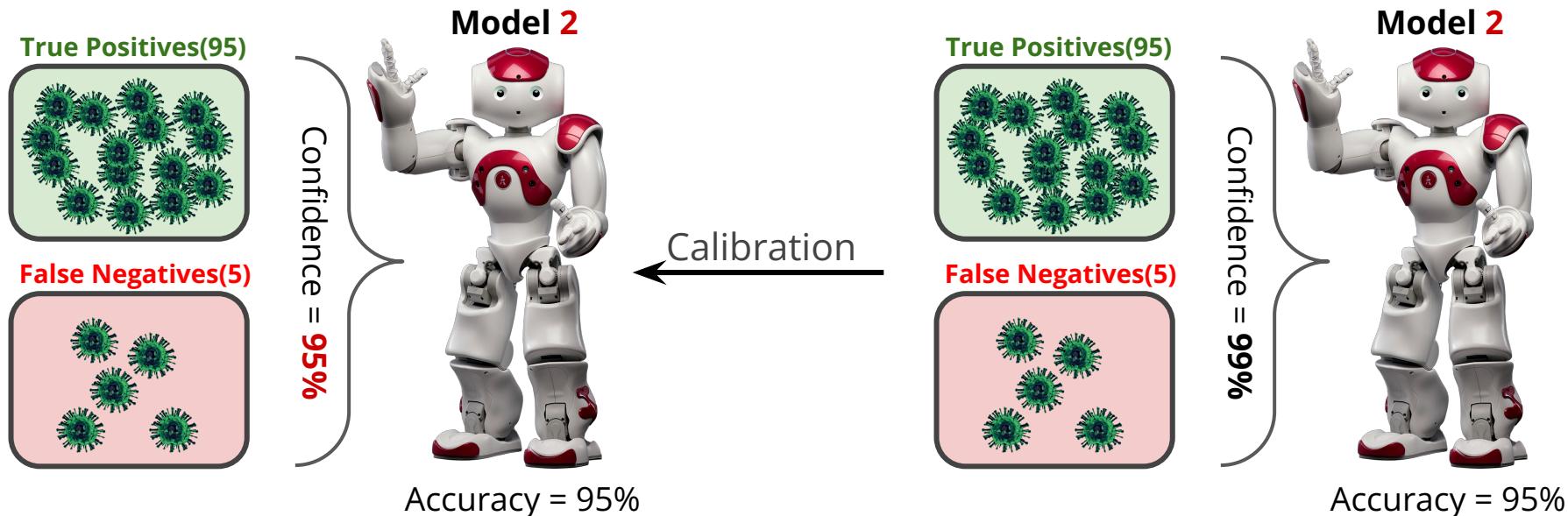
$$290/1000 * |0.98-0.92|$$

Lower ECE is better!

Solution: Calibration Process

Calibration of a classifier

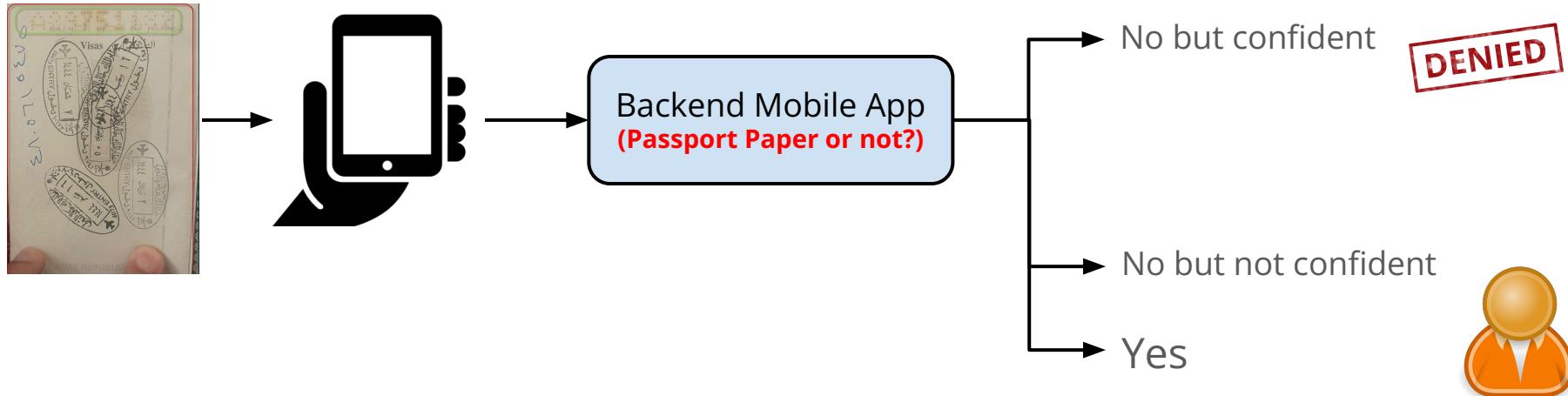
A post-processing step where we take a trained model to improve its predicted confidence levels to match the accuracy of these predictions.



Use Case in Real-Life Scenario

Problem Statement

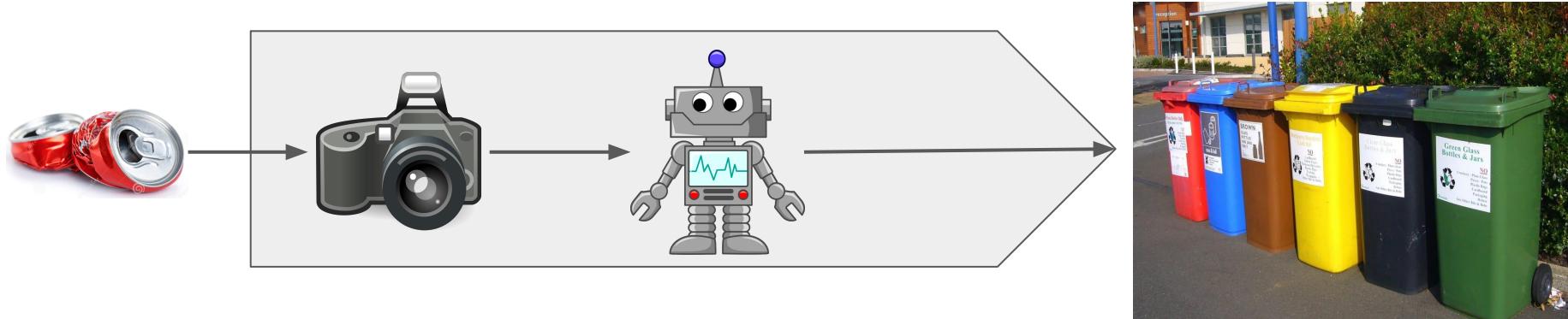
- A customer asked Giza Systems to develop a feature in the customs mobile app, where the user captures an image of his passport with the entry stamp to Egypt.
- A manual effort is done to verify all customs exemption requests.
- Can we reduce efforts by excluding irrelevant images?



Use Case in Real-Life Scenario

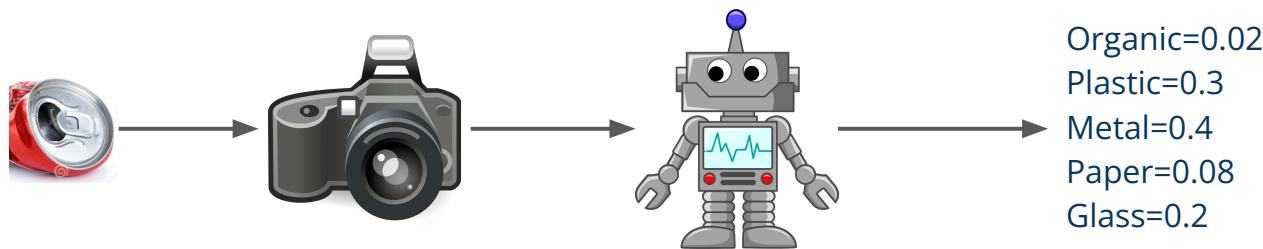
Problem Statement

- A customer asked Giza Systems to develop a smart waste management system. When the user throw a trash in the bin hole, a model classifies its type (Organic, Plastic, Meta, Paper, Glass) and the systems directs it automatically towards the correct waste type pipeline.
- Accurate classification could save much manual effort (time and money)



Use Case in Real-Life Scenario

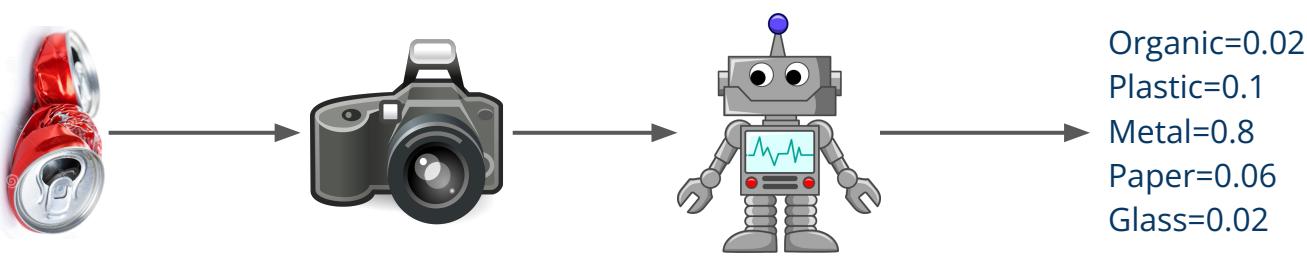
1. Collect dataset for images of trashes of all the needed classes.
2. Train a deep learning model to classify the trashes.
3. Check the calibration performance of your model
 - a. Calibrated → Proceed
 - b. Uncalibrated → Calibrate it OR Re-train with more regularization imposed
4. Is your model confident of the prediction?



40% Confidence = Not Confident

Use Case in Real-Life Scenario

1. Collect dataset for images of trashes of all the needed classes.
2. Train a deep learning model to classify the trashes.
3. Check the calibration performance of your model
 - a. Calibrated → Proceed
 - b. Uncalibrated → Calibrate it OR Re-train with more regularization imposed
4. Is your model confident of the prediction?



80% Confidence = Confident

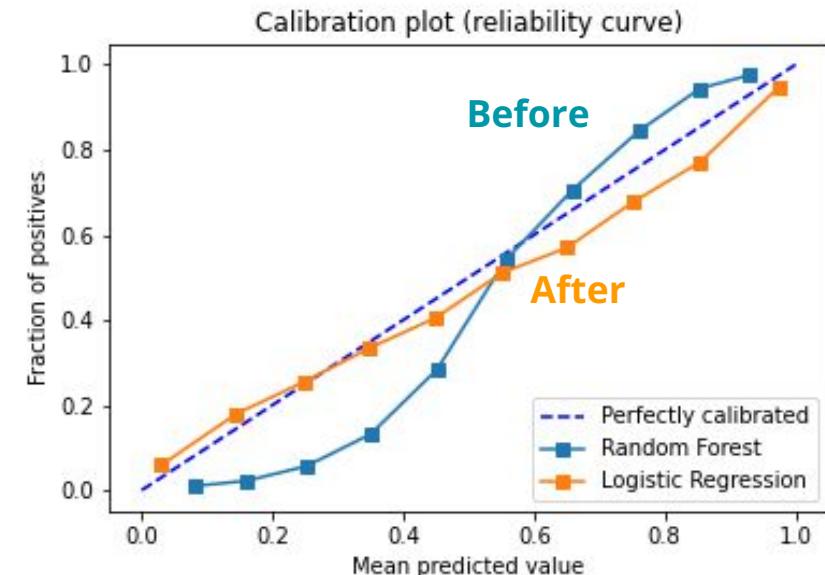
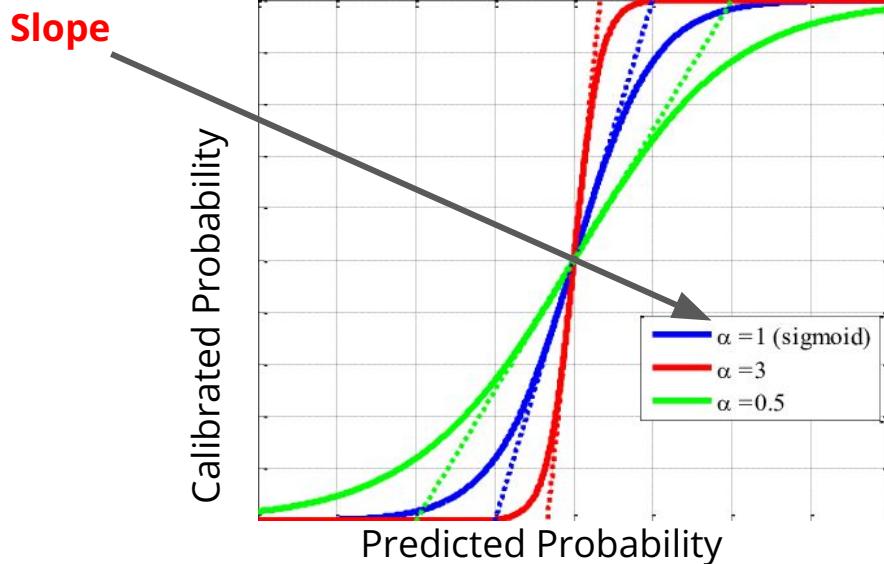
Calibrators

Task	Method	Use Case
Binary Classification	Logistic Calibration [url]	Small datasets
	Beta Calibration [url]	Small to Midsize datasets
	Isotonic Calibration [url]	Large Datasets > 1000 samples to be used in calibration
Multi-Class Classification	Temperature Scaling [url]	Independent of dataset size (Good for Neural Networks)
	Dirichlet Calibration [url]	Needs very large datasets
	Label Smoothing [url]	Occurs during training (No post processing)
	MixUp [url]	Occurs during training (improves generalization and reduce over-confidence)

Examples of Calibrators: Logistic Calibration

Logistic Calibration

Fit the best fitting logistic sigmoid function that maps the model predicted probabilities into calibrated probabilities for a binary class dataset.



Examples of Calibrators: Beta Calibration

Beta Calibration

Fit the best fitting beta function that maps the model predicted probabilities into calibrated probabilities for a binary class dataset.

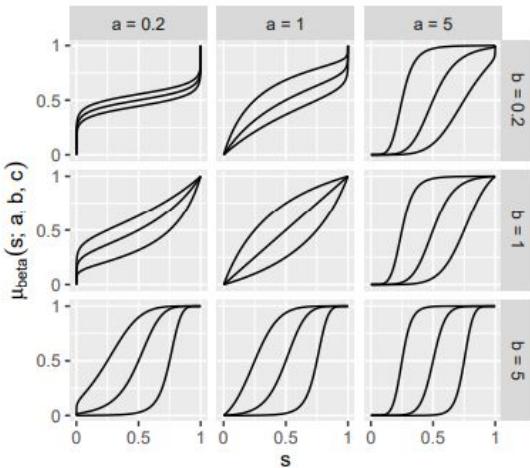
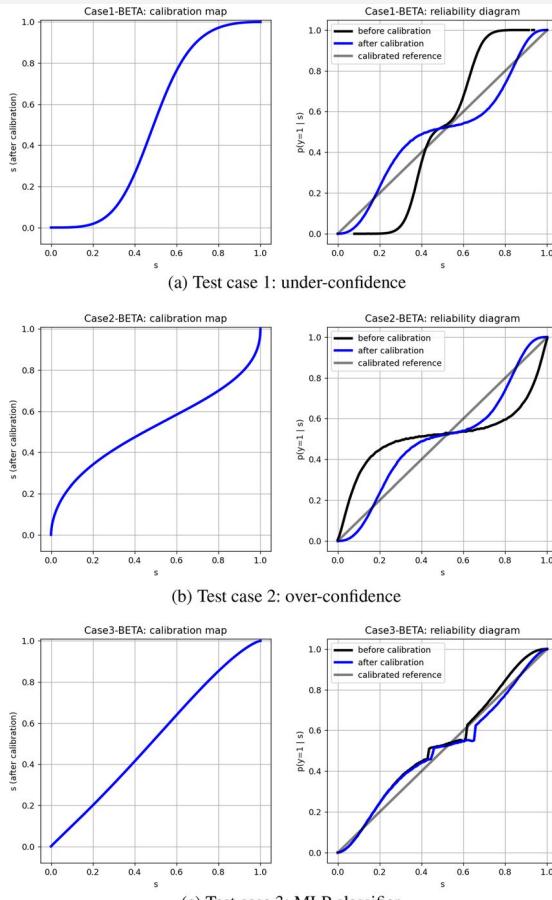


Figure 3: Examples of beta curves with parameters $a, b \in \{0.2, 1, 5\}$, $m \in \{0.25, 0.5, 0.75\}$ and $c = b \ln(1 - m) - a \ln m$.

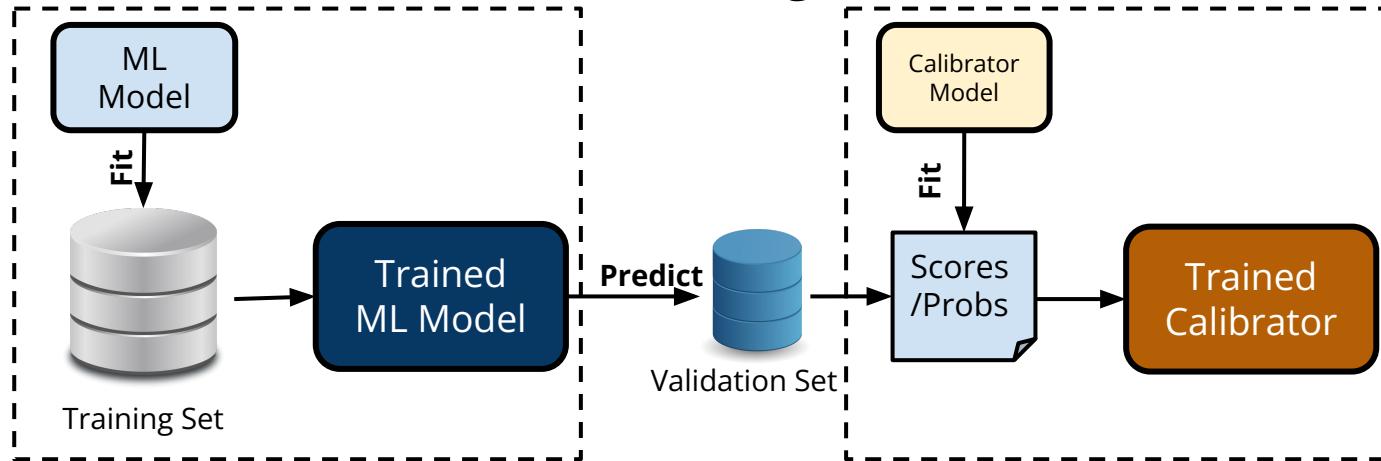
$$\mu_{beta}(s; a, b, c) = \frac{1}{1 + 1 / \left(e^c \frac{s^a}{(1-s)^b} \right)}$$



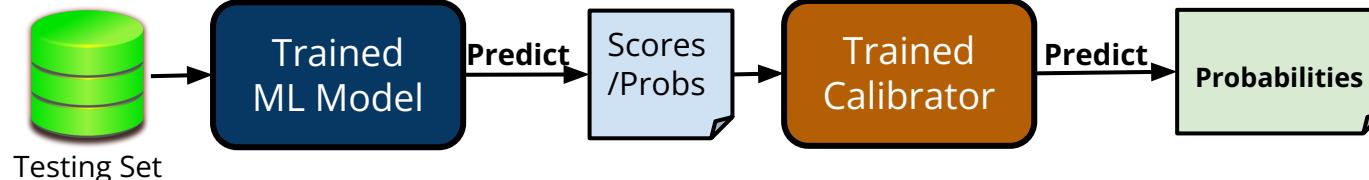
Examples of Calibrators:



Training

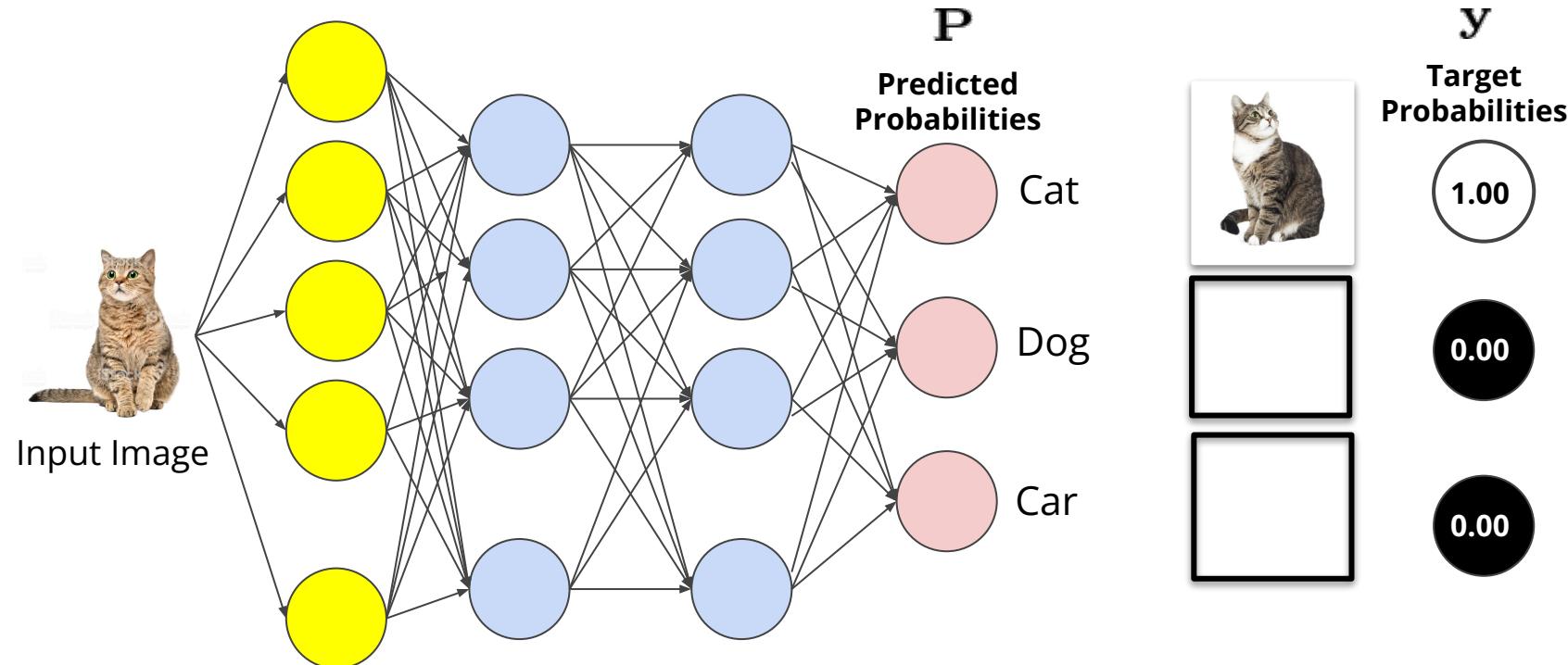


Inference

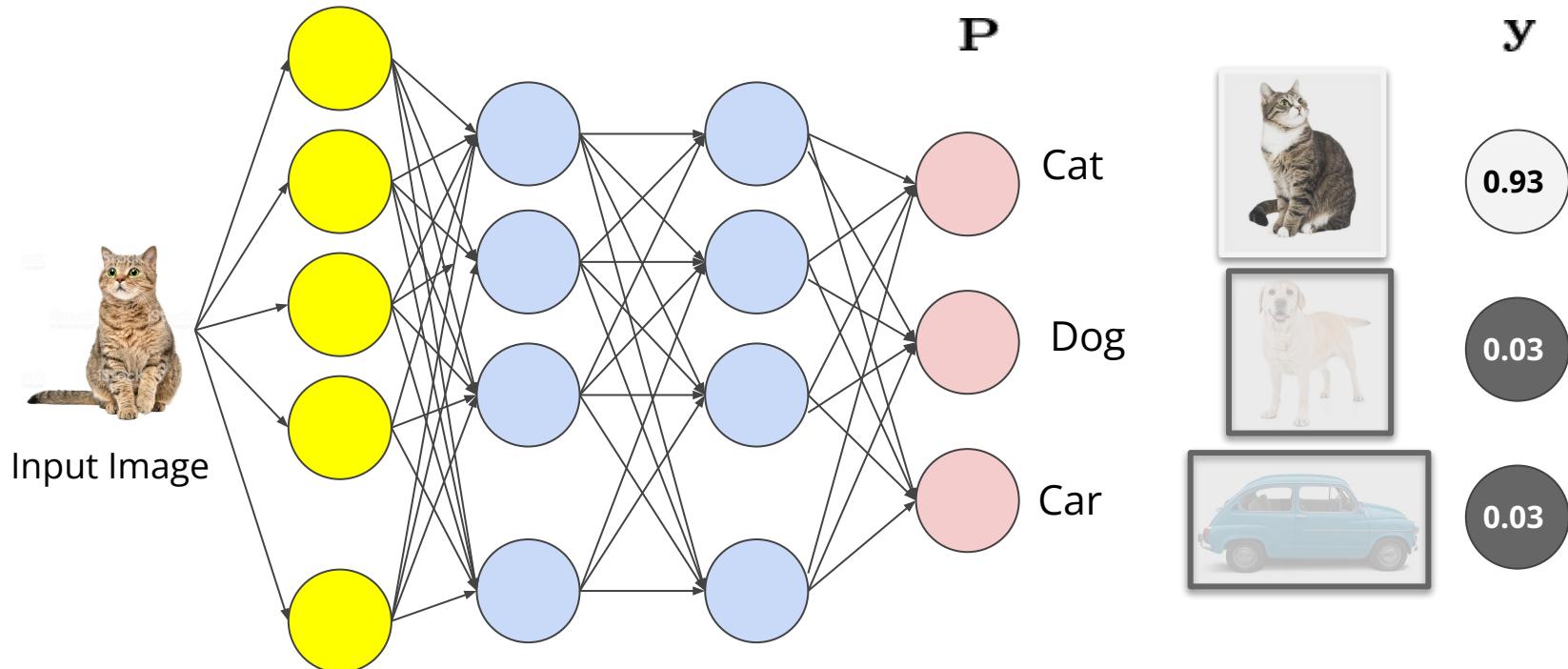


Examples of Calibrators: Label Smoothing

Teach to predict a cat with full confidence.

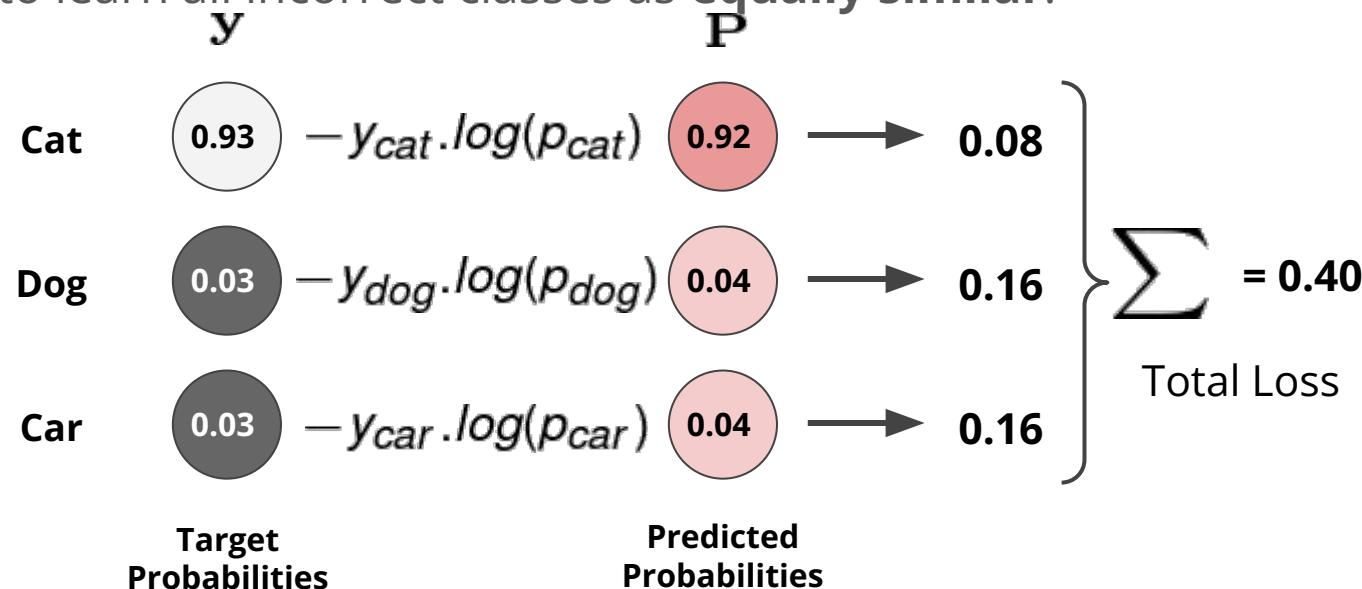


Examples of Calibrators: Label Smoothing

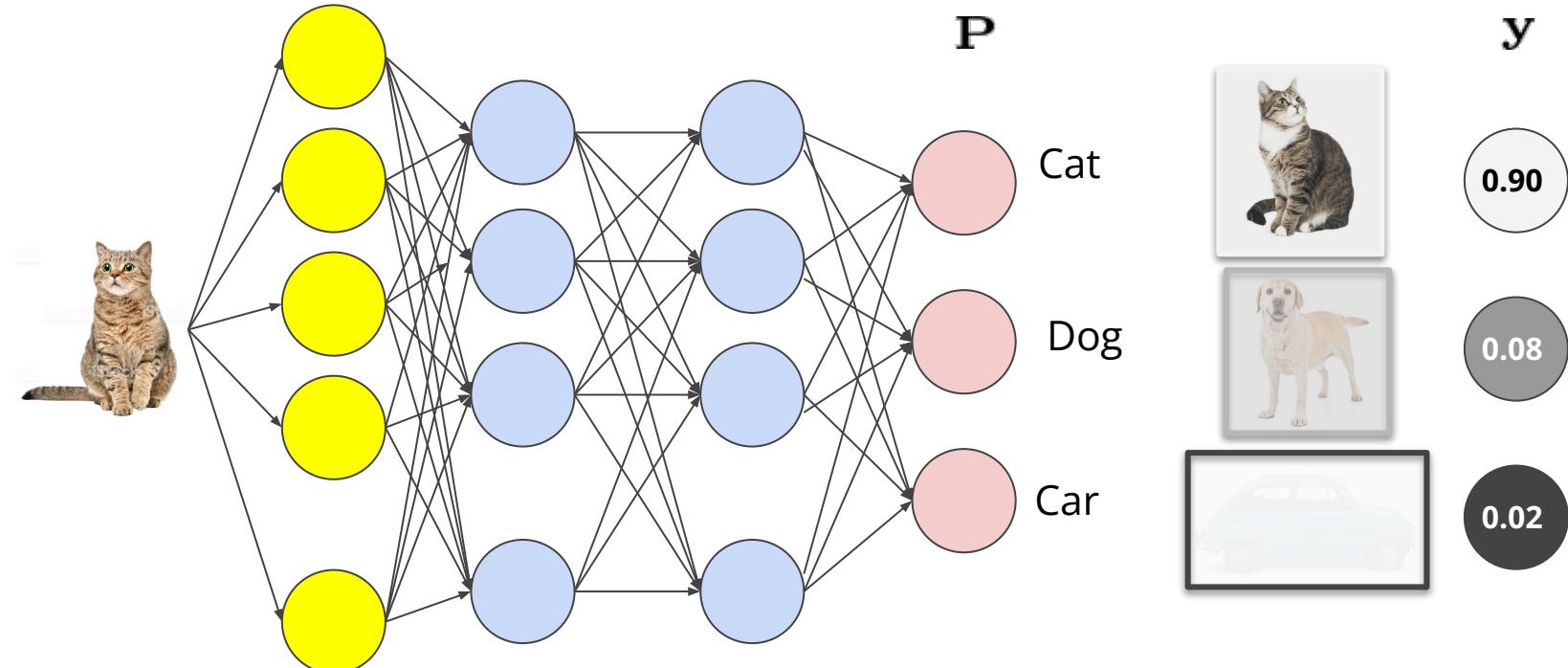


Examples of Calibrators: Label Smoothing

- Less confident predictions
- Enforce to learn all incorrect classes as **equally similar**.



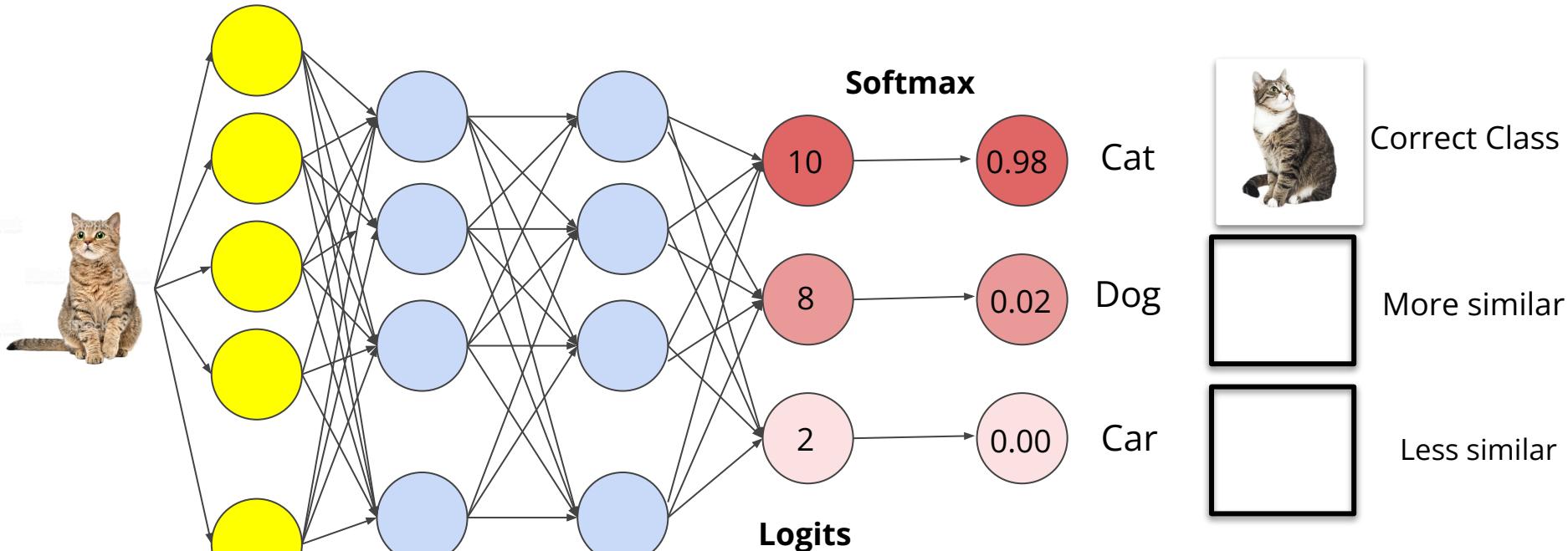
Instance-based label smoothing



M. Maher and M. Kull, "Instance-based Label Smoothing For Better Calibrated Classification Networks," 2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA), Pasadena, CA, USA, 2021, pp. 746-753, doi: 10.1109/ICMLA52953.2021.00124.

Remember: Training without label smoothing

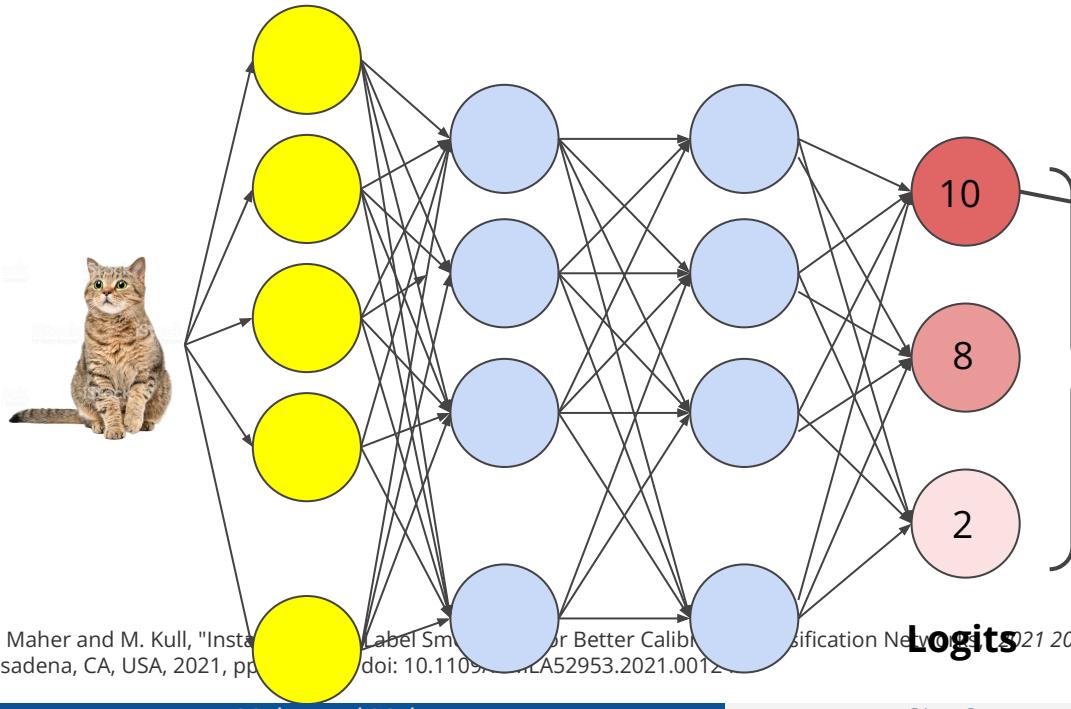
- No enforcement to learn a specific distribution for incorrect classes.



M. Maher and M. Kull, "Instance-Based Label Smoothing For Better Calibrated Classification Networks," 2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA), Pasadena, CA, USA, 2021, pp. 746-753, doi: 10.1109/ICMLA52953.2021.00124.

Choosing the smoothing factor

- Less confident predictions.



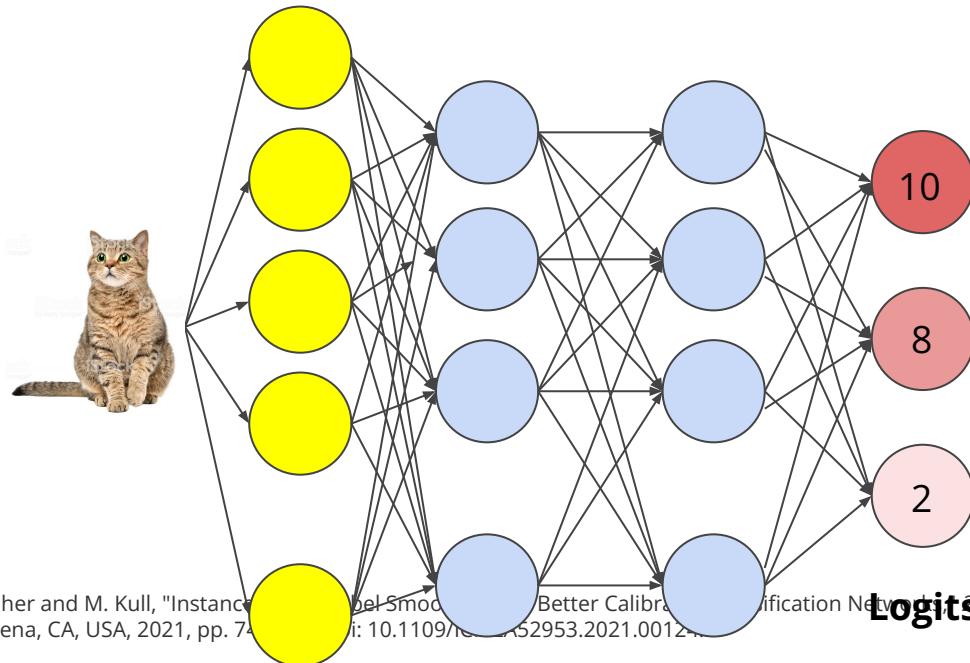
$$\epsilon' = 0.2$$

The more overconfidence,
the higher smoothing factor

$$\epsilon = \frac{10}{10 + 8 + 2} * \epsilon' = 0.1$$

Instance-based label smoothing

- Less confident predictions.
- Enforce to learn classes' similarity of the original dataset.



$$\epsilon = 0.1$$

$$y_{cat} = 1 - \epsilon = 0.9$$

0.90

$$y_{dog} = \frac{8}{8 + 2} \epsilon = 0.08$$

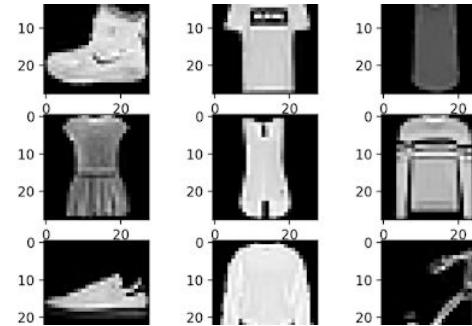
0.08

$$y_{car} = \frac{2}{8 + 2} \epsilon = 0.02$$

0.02

Experimental Results

Fashion-MNIST



Cifar-10



Cifar-100

Experimental Results

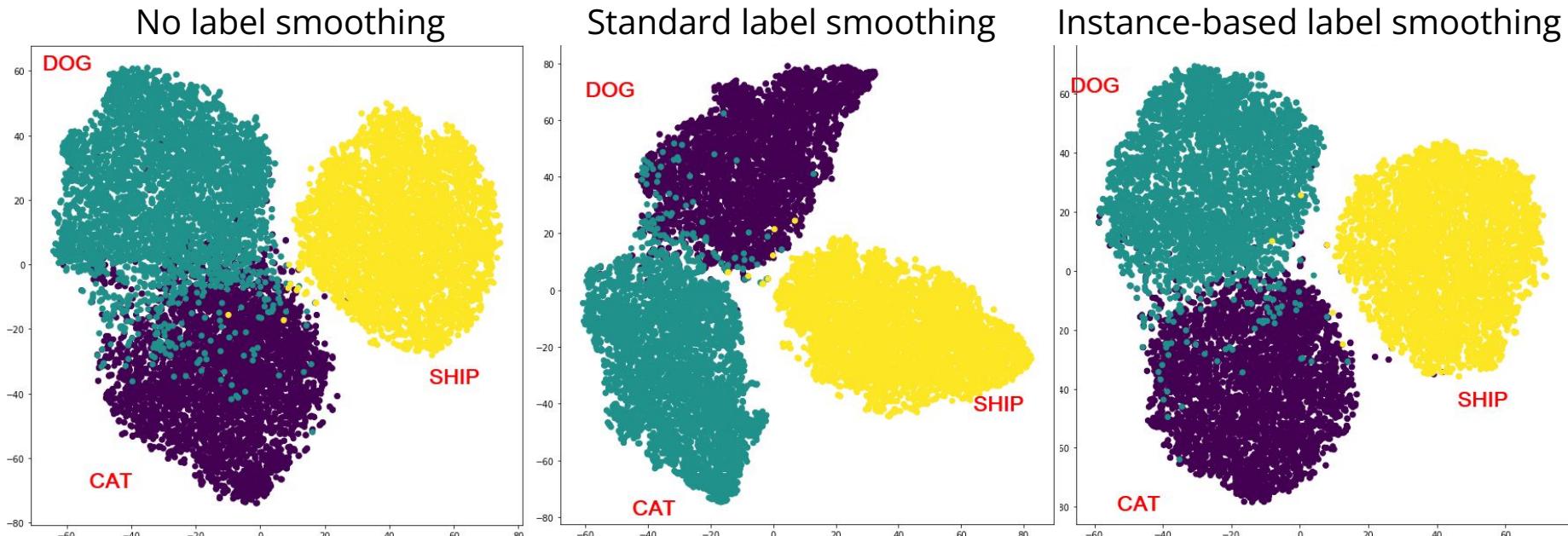


Figure: T-SNE visualization of the output logits for 3 different classes of Cifar10 dataset.

Temperature Scaling

- A 3-class dataset (cats, dogs and cars)

Logits: [10.0, 6.0, 1.0]

Softmax

Probabilities: [0.982, 0.018, 0.0001]

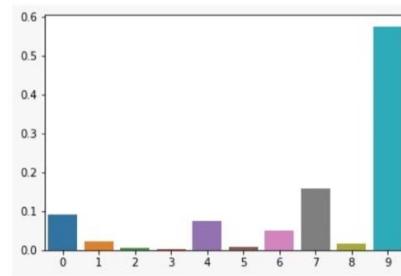
Logits: [10/3, 6/3, 1/3]
Temperature: 3

Temp Scaling +
Softmax

Probabilities: [0.761, 0.201, 0.038]

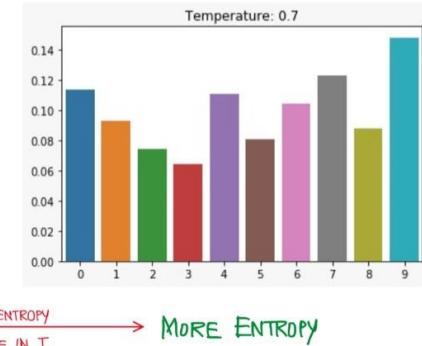
SOFTMAX WITHOUT TEMPERATURE ($T=1$)

$$\frac{e^{z_i}}{\sum_j e^{z_j}}$$



SOFTMAX WITH TEMPERATURE

$$\frac{e^{z_i/\tau}}{\sum_j e^{z_j/\tau}}$$



Temperature Scaling

- A 3-class dataset (cats, dogs and cars)

Logits: [10.0, 6.0, 1.0]

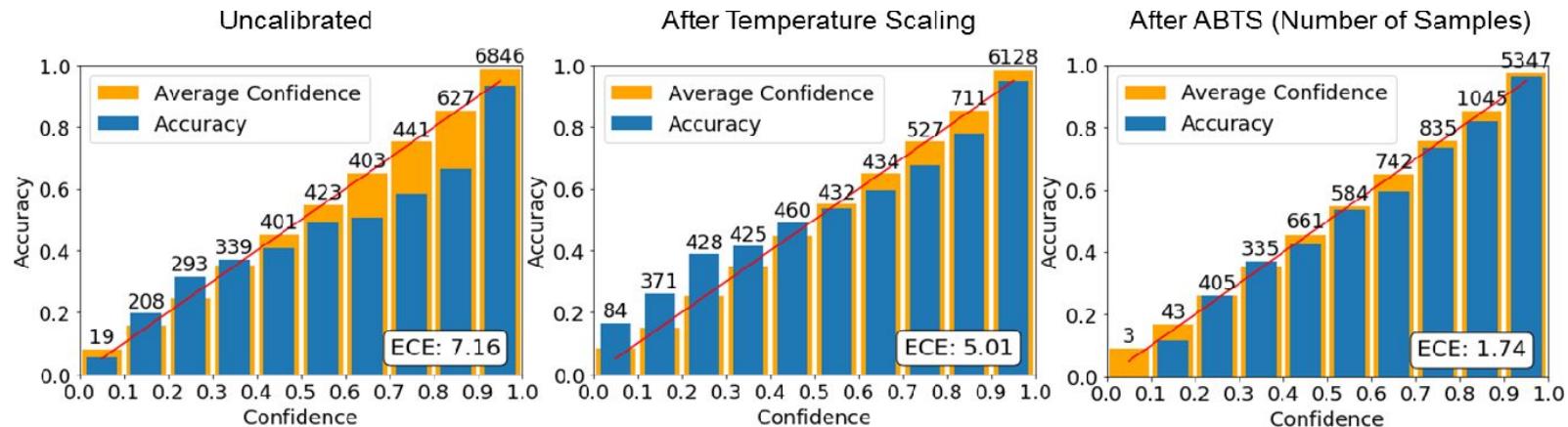
Softmax

Probabilities: [0.982, 0.018, 0.0001]

Logits: [10/3, 6/3, 1/3]
Temperature: 3

Temp Scaling +
Softmax

Probabilities: [0.761, 0.201, 0.038]



MixUp: Reduce Neural Networks Over-confidence

1. **Data Interpolation:** Given two training examples $(\mathbf{x}_i, \mathbf{y}_i)$ and $(\mathbf{x}_j, \mathbf{y}_j)$, mixup generates a new training example $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$

$$\tilde{\mathbf{x}} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j$$

$$\tilde{\mathbf{y}} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$$

where λ is a mixing coefficient sampled from a Beta distribution, typically $\lambda \sim \text{Beta}(\alpha, \alpha)$ for some parameter $\alpha > 0$.
A common choice is $\alpha=0.2$.

2. **Label Interpolation:** The labels are also linearly interpolated, so instead of assigning a hard class label, the model learns from a *soft label*, representing a mix of the two classes.

For instance, if \mathbf{y}_i is the label for a cat image (e.g., one-hot encoded as [1,0]) and \mathbf{y}_j is the label for a dog image [0,1], then the synthetic label might be something like [0.8,0.3], indicating 80% "cat" and 20% "dog".



Zhang, Hongyi. "mixup: Beyond empirical risk minimization." *arXiv preprint arXiv:1710.09412* (2017). → Published in ICLR 2018

Takeaways

- Do not get fooled with “highly accurate models”.
 - Metrics like Accuracy, F-Score, Precision, Recall cares only about the class label
- Probabilistic Classification is important in cost sensitive domains.
 - I can discard unconfident model predictions.
 - I can use another model to improve the prediction confidence.
- Machine Learning Classifiers can be used to predict scores/probabilities
- These probabilities are usually uncalibrated and need further calibration methods.
- Calibration performance is measured through Expected Calibration Error (ECE) and Reliability diagram.
- Calibration can be done using methods like logistic calibration and isotonic calibration for binary classification, or temperature scaling and label smoothing for multi-class classification.

Takeaways

- Can I use binary classification calibration methods in multi-classification task?
 - **Yes, 1 Vs All**
- Good Research Point:
 - Automation of the Calibration / Over-confidence Reduction?

Measuring Uncertainty in Regression Models: Bootstrapping

- Jack-Knife+:
 - **leave-one-out** method consists of training a model for each data point in our dataset, training it on the entire dataset but removing one sample at a time.

F1	F2	F3	Tar
..
..
..
..
..
..
..

Model 1

Kim, Byol, Chen Xu, and Rina Barber. "Predictive inference is free with the jackknife+-after-bootstrap." *Advances in Neural Information Processing Systems* 33 (2020): 4138-4149.

Measuring Uncertainty in Regression Models: Bootstrapping

- Jack-Knife+:
 - **leave-one-out** method consists of training a model for each data point in our dataset, training it on the entire dataset but removing one sample at a time.

F1	F2	F3	Tar	Model 1
..	Model 1
..	Model 2
..	
..	
..	
..	
..	
..	
..	

Kim, Byol, Chen Xu, and Rina Barber. "Predictive inference is free with the jackknife+-after-bootstrap." *Advances in Neural Information Processing Systems* 33 (2020): 4138-4149.

Measuring Uncertainty in Regression Models: Bootstrapping

- Jack-Knife+:
 - **leave-one-out** method consists of training a model for each data point in our dataset, training it on the entire dataset but removing one sample at a time.

F1	F2	F3	Tar	
..	Model 1
..	Model 2
..	Model 3
..	
..	
..	
..	

Kim, Byol, Chen Xu, and Rina Barber. "Predictive inference is free with the jackknife+-after-bootstrap." *Advances in Neural Information Processing Systems* 33 (2020): 4138-4149.

Measuring Uncertainty in Regression Models: Bootstrapping

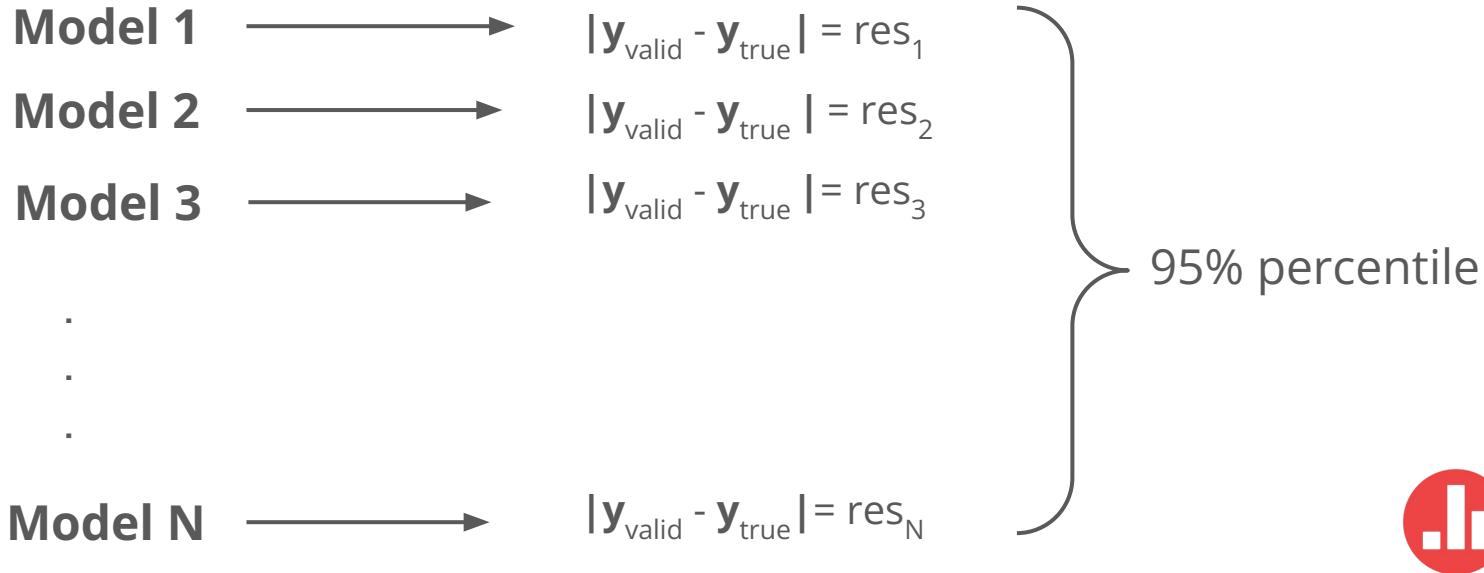
- Jack-Knife++:
 - **leave-one-out** method consists of training a model for each data point in our dataset, training it on the entire dataset but removing one sample at a time.

F1	F2	F3	Tar	
..	Model 1
..	Model 2
..	Model 3
..	
..	
..	
..	
..	Model N
..	

Kim, Byol, Chen Xu, and Rina Barber. "Predictive inference is free with the jackknife+-after-bootstrap." *Advances in Neural Information Processing Systems* 33 (2020): 4138-4149.

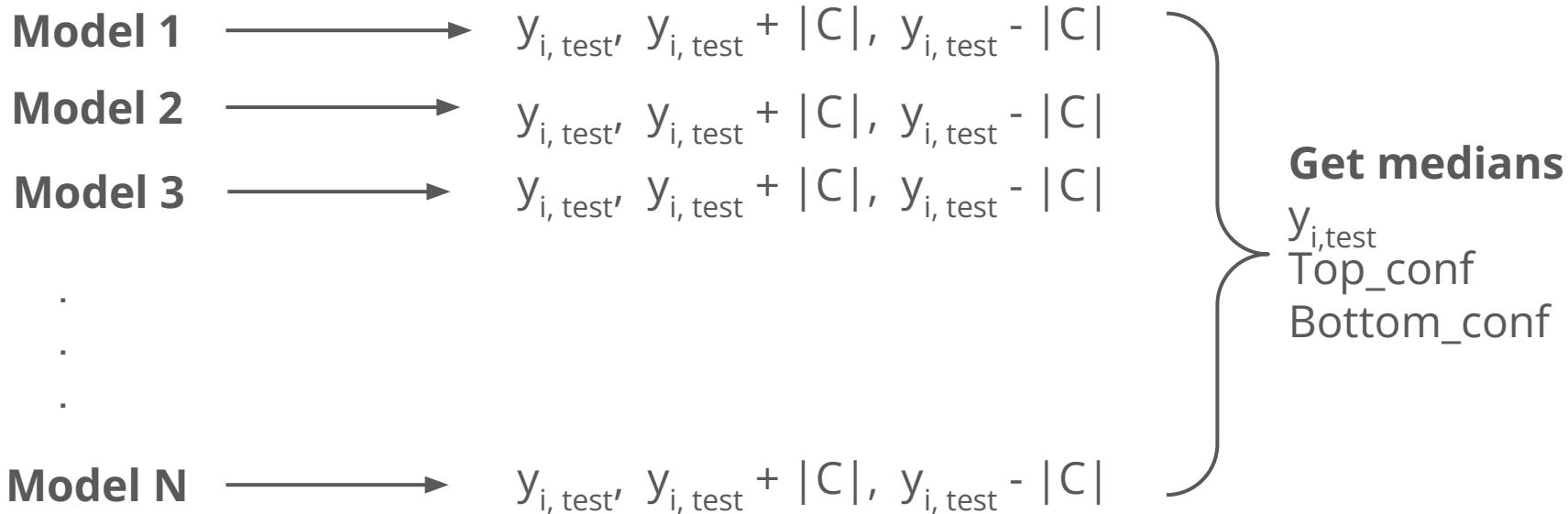
Measuring Uncertainty in Regression Models: Bootstrapping

- Jack-Knife+:
 - Get residuals from each model on the validation set.
 - Compute the confidence interval thresholds for the N-percentile.

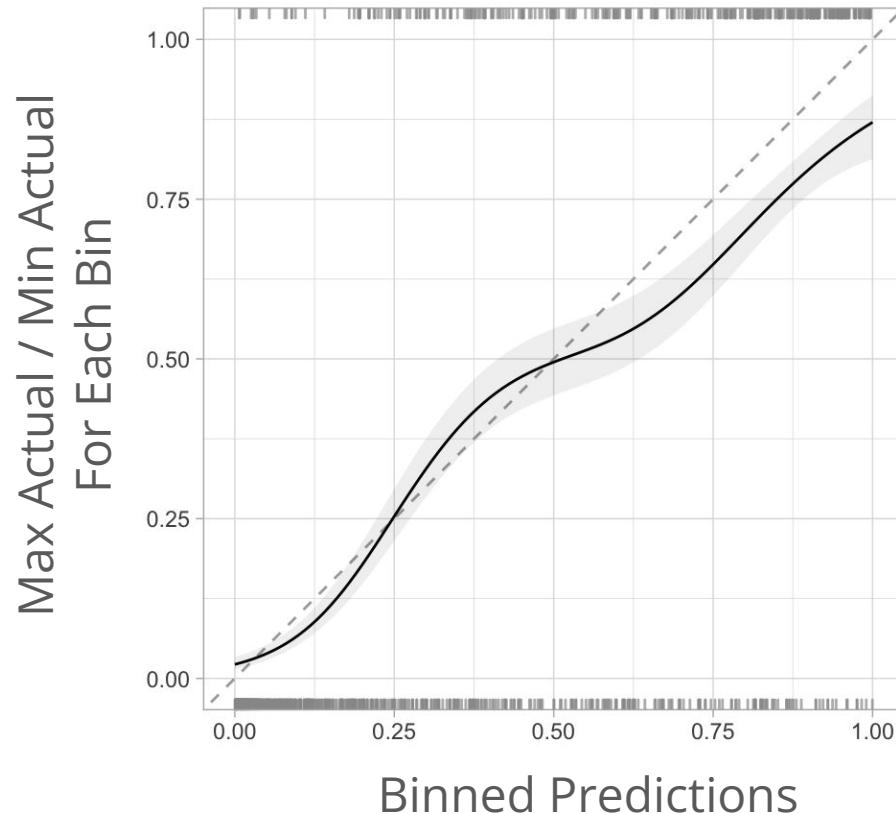


Measuring Uncertainty in Regression Models: Bootstrapping

- Jack-Knife+:
 - Compute the prediction on the test set instance
 - Add the confidence interval threshold



Calibration Plot for Regression Models



Hands-On Task (New Task)

- **Data Team Details:**

- Individual Submissions
- Due date: **11 November (Data)**
- Delivery via email to mohamed.maher@gizasystems.com & headway 2024 data testing@gizasystems.com
 - Title: "**Calibration Task Delivery**"
 - Content:
 - Model Training Code files (notebook/.py),
 - Describe briefly your approach and selected models,
 - Pipeline evaluation file (.py) to be deployed as an endpoint when executed.
 - Endpoint
 - Body:
 - {0: {'meta-feature1_name': value, 'meta-feature2_name': value,},
1: {'meta-feature1_name': value, 'meta-feature2_name': value,}, ...}
 - Response:
 - {0: {'XGBRegressor': 0.73, 'LinearSVR': 0.14,},
1: {'LASSO': 0.63, 'LinearSVR': 0.26, ...}}

Hands-On Task (New Task)

- **Data Testing Team Details**

- Work in pairs
- Due date: **17 November (Data Testing)**
- Delivery via email to mohamed.maher@gizasystems.com
 - Title: "**Calibration Task Evaluation Delivery**"
 - Content:
 - The script used to evaluate the models
 - Report
 - Team members and who worked on which models?
 - Comparison:
 - Measuring the accuracy (doesn't mean accuracy metric) performance of the models.
 - F1, AUC, Micro Recall, ...?
 - Justify your selection
 - Measuring the reliability of the models.
 - Reliability diagram
 - ECE