

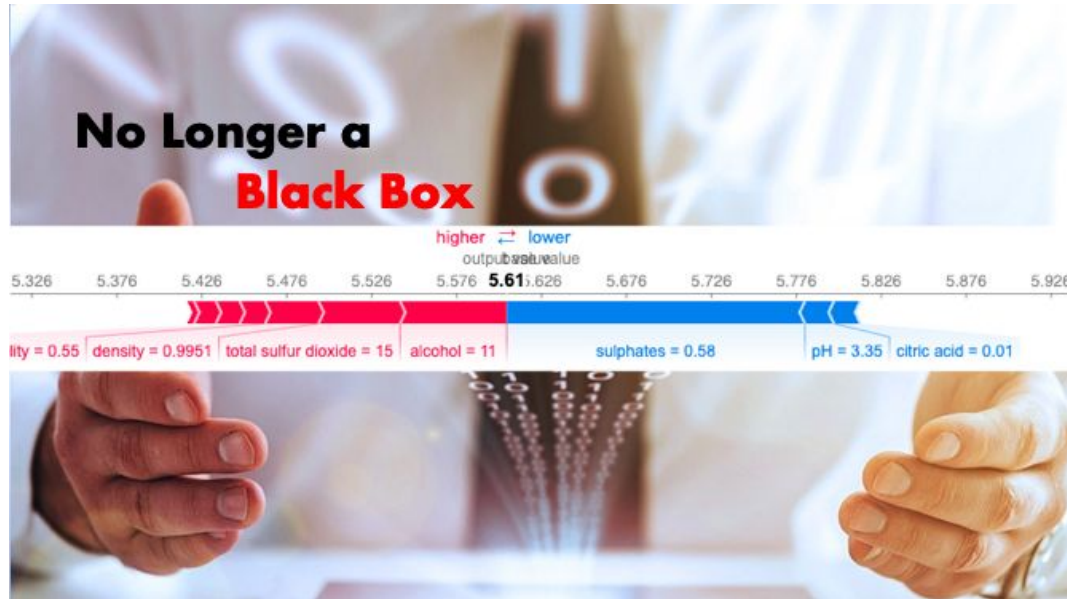
A Unified Approach to Interpreting Model Predictions

Teerapong Panboonyuen, Ph.D.

<https://kaopanboonyuen.github.io/>

Reference

Lundberg, Scott M., and Su-In Lee.
"A unified approach to interpreting model predictions."
Advances in neural information processing systems. 2017.



Lundberg and Lee (NIPS 2017)

SHAP **SHAPLEY ADDITIVE** **EXPLANATIONS**

Lundberg and Lee (NIPS 2017) - Main (1)

They proposed the **SHAP value** as a united approach to explaining the output of any machine learning model.

Lundberg and Lee (NIPS 2017) - Main (1)

They proposed the **SHAP value** as a united approach to explaining the output of any machine learning model.

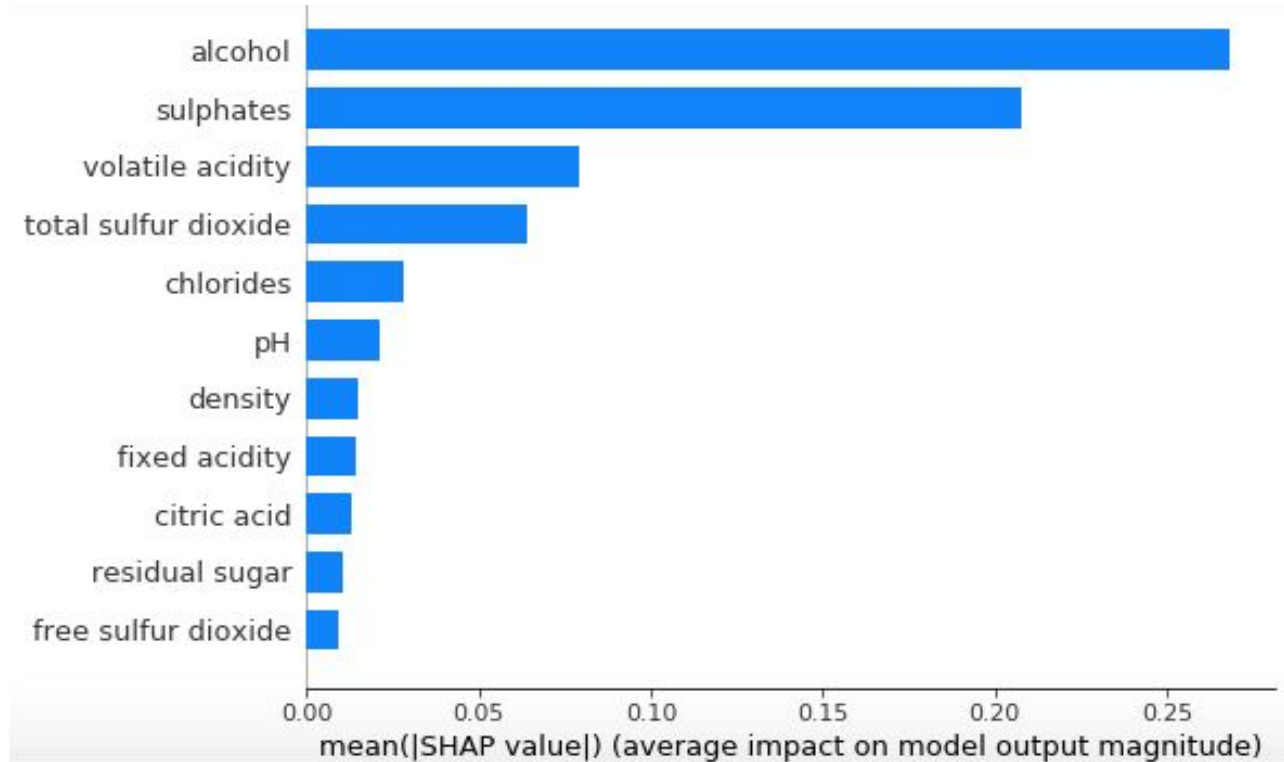
Three benefits worth mentioning here.

The first one is **global interpretability**

— the collective SHAP values can show how much each predictor contributes, either positively or negatively, to the target variable.

- (a) This is like the variable importance plot but it is able to show the positive or negative relationship for each variable with the target

Variable Importance Plot



Lundberg and Lee (NIPS 2017) - Main (1)

They proposed the **SHAP value** as a united approach to explaining the output of any machine learning model.

The second benefit is ***local interpretability*** — each observation gets its own set of SHAP values. This greatly increases its transparency.

- (a) We can explain why a case receives its prediction and the contributions of the predictors. Traditional variable importance algorithms only show the results across the entire population but not on each individual case.
- (b) The local interpretability enables us to pinpoint and contrast the impacts of the factors.

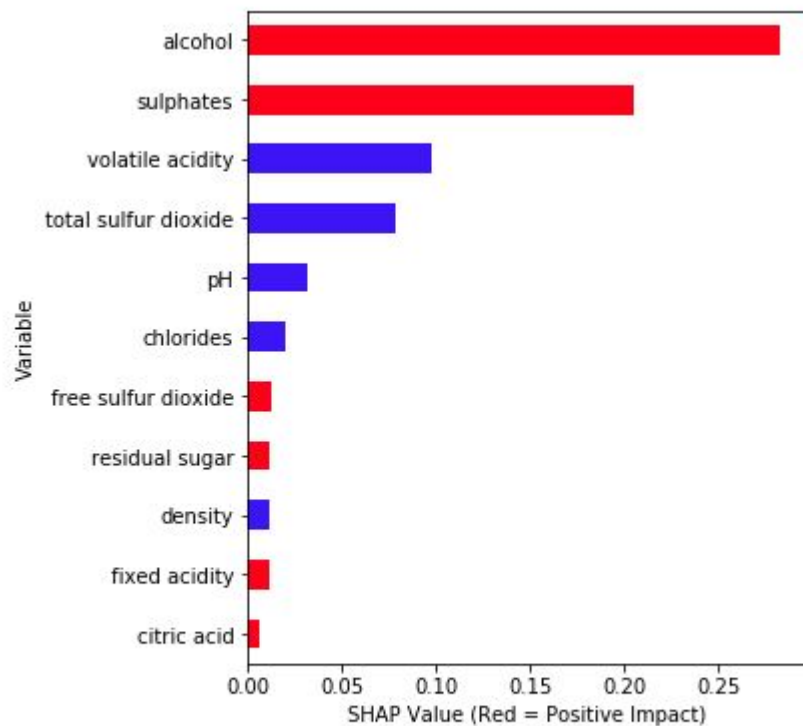
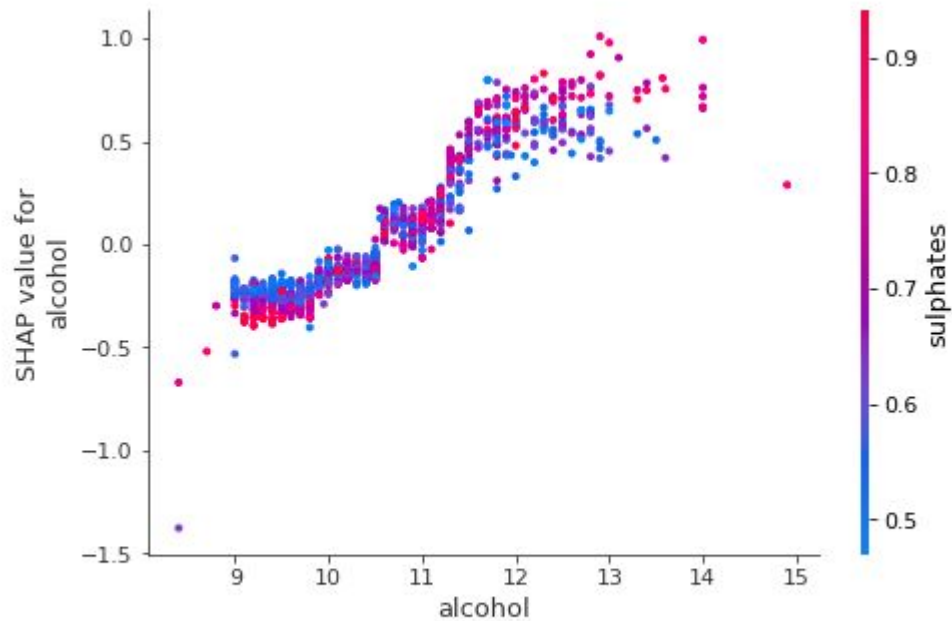


Exhibit (K.1): The simplified version



The SHAP Dependence Plot

N O T E

SHAP values can be calculated for **any tree-based model**,
while other methods use linear regression or logistic
regression models
as the ***surrogate models***.

Chris Kuo, Ph.D. | Columbia University

Model Interpretability Does Not Mean Causality.

What is the Shapley Value?



Let me explain the Shapley value with a story

Assume Ann, Bob, and Cindy together were hammering an “error” wood log, 38 inches, to the ground.

After work, they went to a local bar for a drink and I, a mathematician, came to join them.

I asked a very bizarre question:
“What is everyone’s contribution (in inches)?”

That's the way to calculate the Shapley value

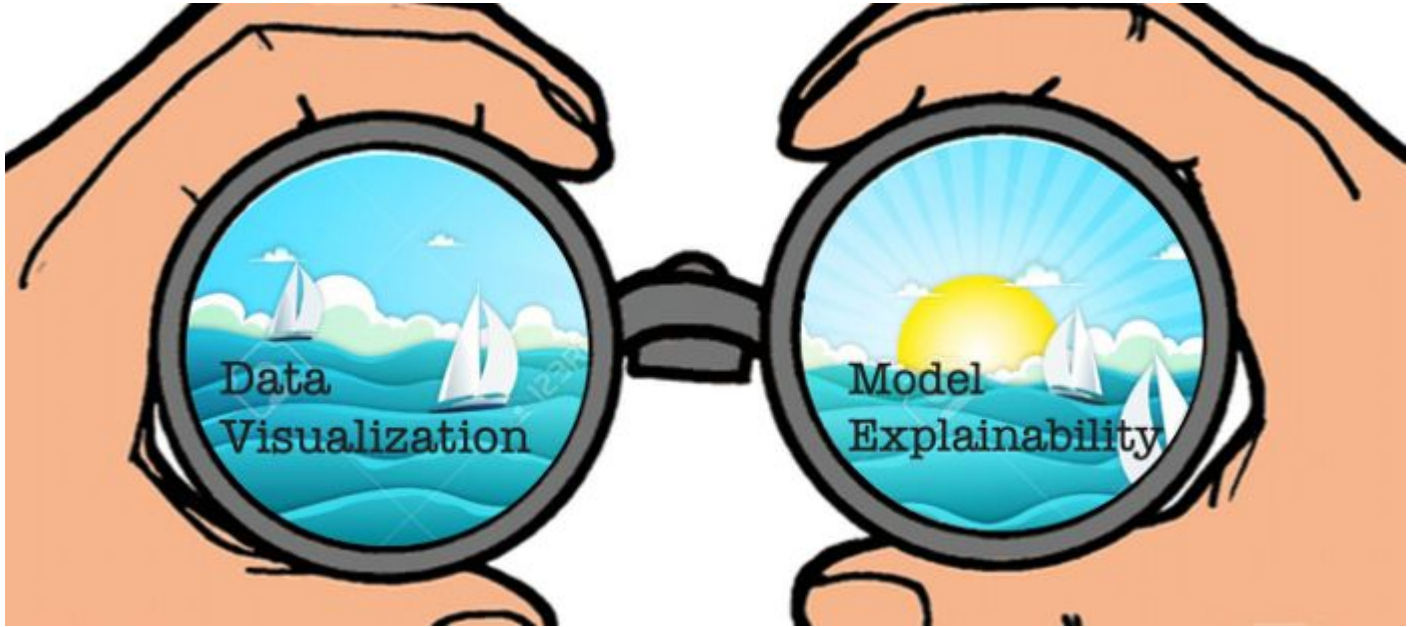
It is the average of the marginal contributions across all permutations.

How do we **measure the contributions of the hammers** (predictors)?

The Shapley values!

Combination	Marginal contribution			inches
	Ann	Bob	Cindy	Total
A, B, C	2	32	4	38
A, C, B	4	34	0	38
B, A, C	2	32	4	38
B, C, A	0	28	10	38
C, A, B	2	36	0	38
C, B, A	0	28	10	38
Average	2	32	4	38

Data Visualization and Model Explainability



Example: with my code

(A) Variable Importance Plot — Global Interpretability

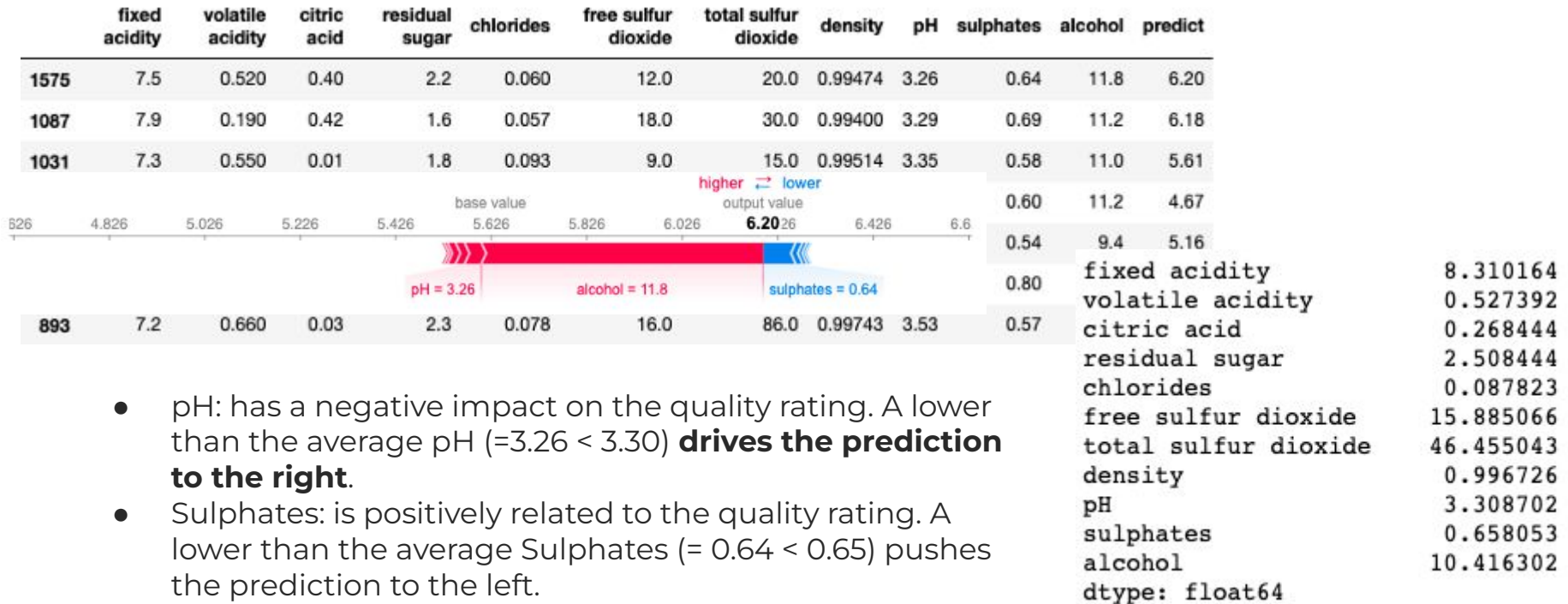
(B) SHAP Dependence Plot — Global Interpretability

(C) Individual SHAP Value Plot — Local Interpretability

Alcohol: has a positive impact on the quality rating. The alcohol content of this wine is 11.8 (as shown in the first row of Table B) which is higher than the average value 10.41. **So it pushes the prediction to the right.**



Example: with my code



- pH: has a negative impact on the quality rating. A lower than the average pH ($=3.26 < 3.30$) **drives the prediction to the right.**
- Sulphates: is positively related to the quality rating. A lower than the average Sulphates ($= 0.64 < 0.65$) pushes the prediction to the left.

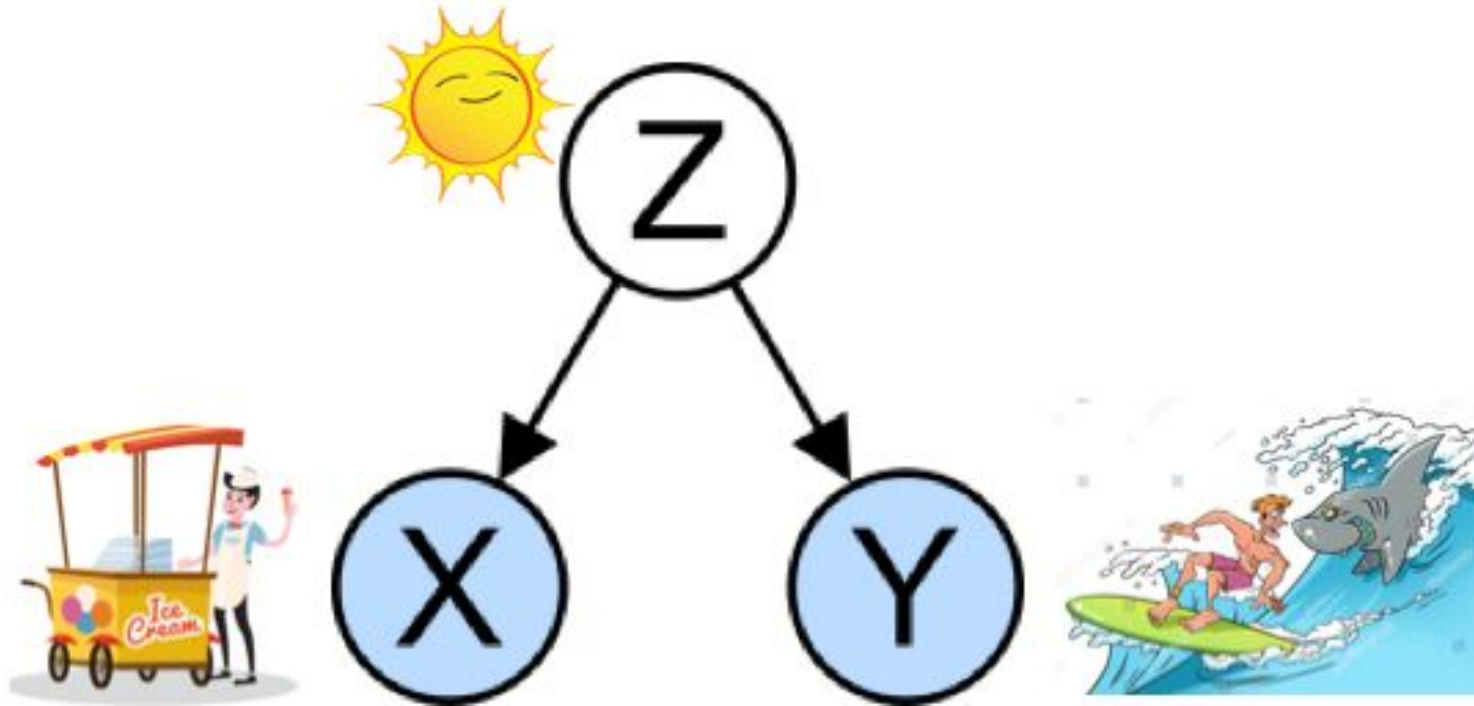
Correlation Can Mean Causation — Only When Certain Conditions Meet



Inferring Causation from Correlation Is a Scary Thing



Confounding Factor



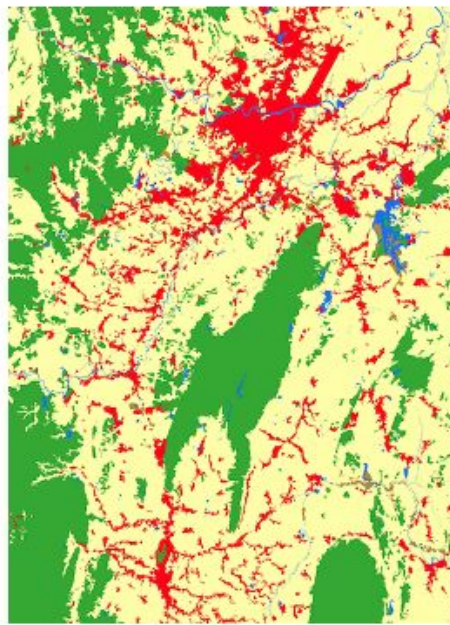
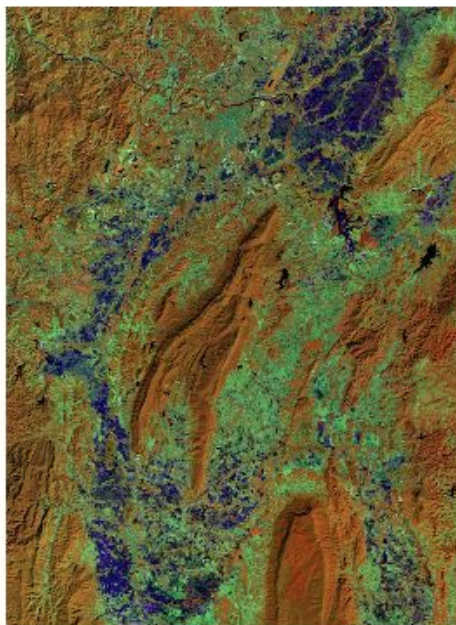
Efficient Uncertainty Estimation for Semantic Segmentation

Teerapong Panboonyuen, Ph.D.

<https://kaopanboonyuen.github.io/>

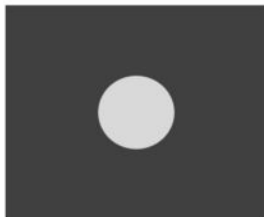
Public and Private Corpora

Private corpus (GISTDA Nan Province Corpus)

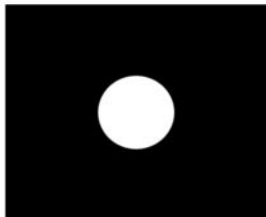


Color	Class
Yellow	Agriculture
Green	Forest
Brown	Miscellaneous
Red	Urban
Blue	Water

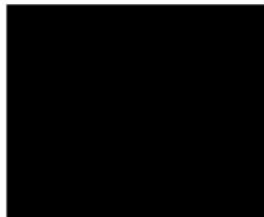
Input image of
aircraft part



Neural network
segmentation

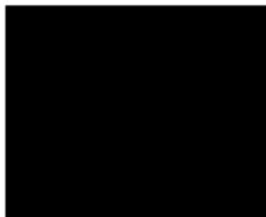


Uncertainty map
from BCNN

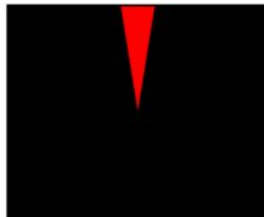
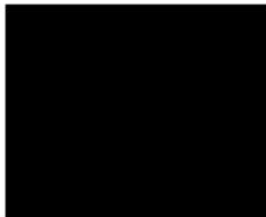


Production decision
based on segmentation

Neural network successfully
recognizes circle as void in material.
Defective with *low uncertainty*.

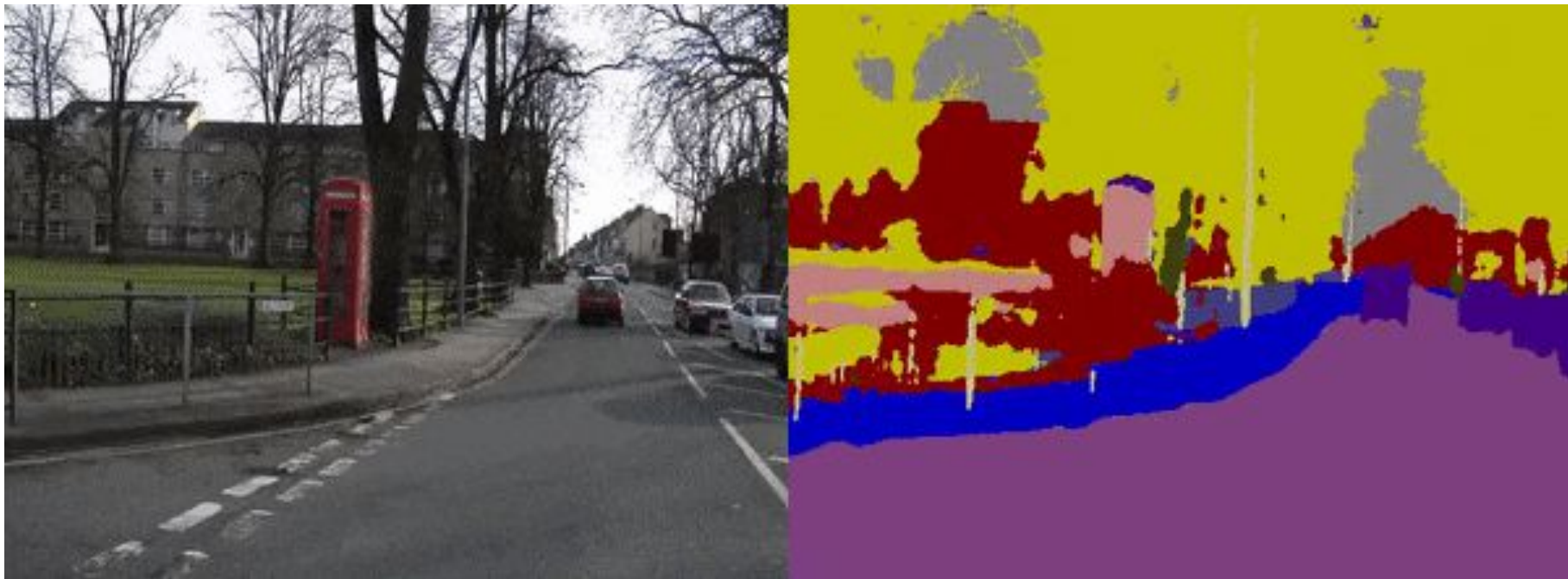


Neural network fails to recognize low-
contrast circle as void in material.
Non-defective with *high uncertainty*.



Neural network was not trained on
parts with cracks, so it fails to
recognize triangle as crack in material.
Non-defective with *high uncertainty*.

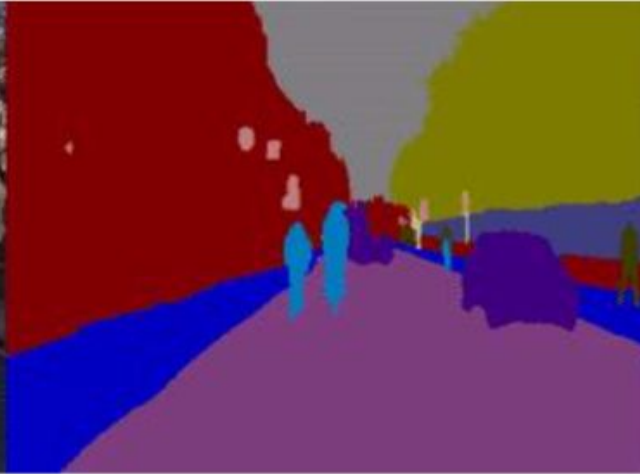
Predictions from Tiramisu on CamVid video stream.



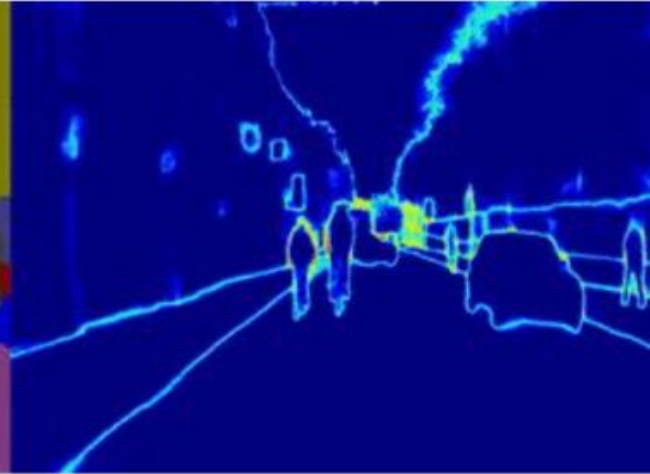
Bayesian SegNet for probabilistic scene understanding



Input Image



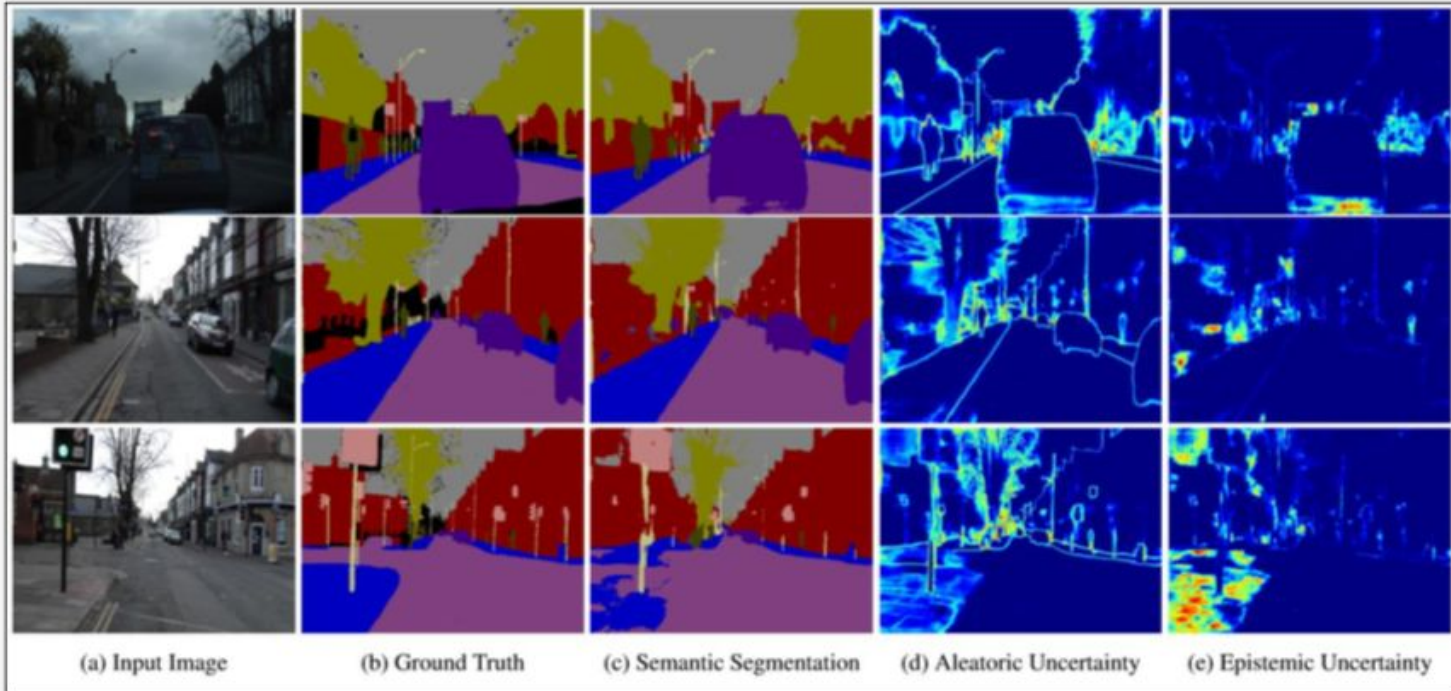
Semantic
Segmentation



Uncertainty

What kind of uncertainty can we model?

Epistemic uncertainty is modeling uncertainty
Aleatoric uncertainty is sensing uncertainty



Modeling Uncertainty with Bayesian Deep Learning



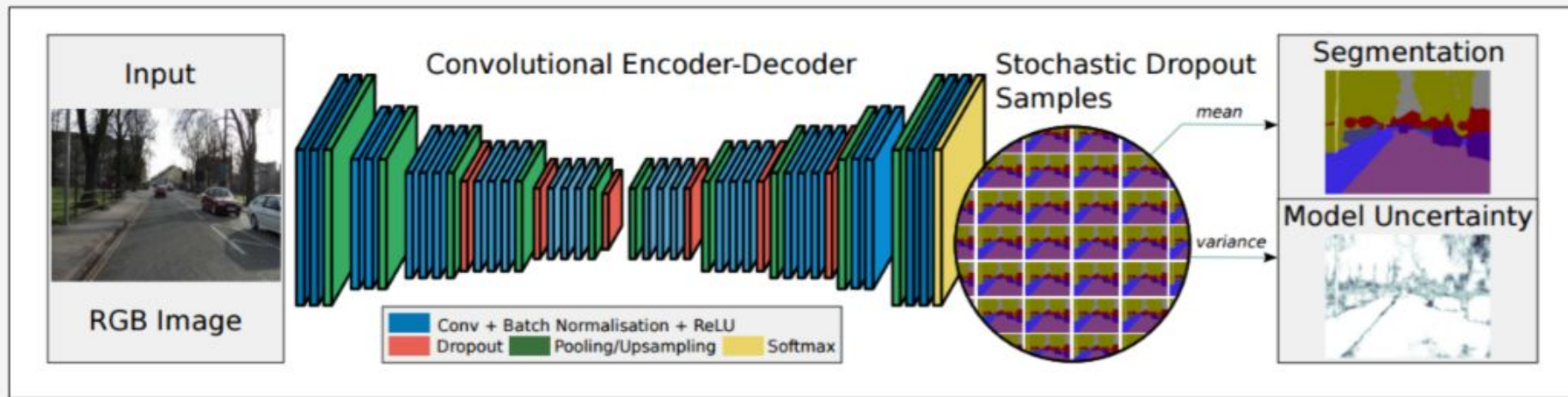
Deep learning is required to achieve state of the art results in computer vision applications but doesn't provide uncertainty estimates.

- **Bayesian neural networks** are a framework for understanding uncertainty in deep learning
- They have **distributions over network parameters** (rather than deterministic weights)
- Traditionally they have been **tricky to scale**

Modeling Epistemic Uncertainty with Bayesian Deep Learning

We can **model epistemic uncertainty** in deep learning models using Monte Carlo **dropout sampling** at test time.

Dropout sampling can be interpreted as **sampling from a distribution over models**.



Modeling Aleatoric Uncertainty with Probabilistic Deep Learning

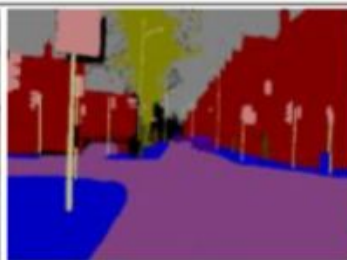
	Deep Learning	Probabilistic Deep Learning
Model	$[\hat{y}] = f(x)$	$[\hat{y}, \hat{\sigma}^2] = f(x)$
Regression	$Loss = \ y - \hat{y}\ ^2$	$Loss = \frac{\ y - \hat{y}\ ^2}{2\hat{\sigma}^2} + \log \hat{\sigma}^2$
Classification	$Loss = SoftmaxCrossEntropy(\hat{y}_t)$	$\hat{y}_t = \hat{y} + \epsilon_t \quad \epsilon_t \sim N(0, \hat{\sigma}^2)$ $Loss = \frac{1}{T} \sum_t SoftmaxCrossEntropy(\hat{y}_t)$

Semantic Segmentation Performance on CamVid

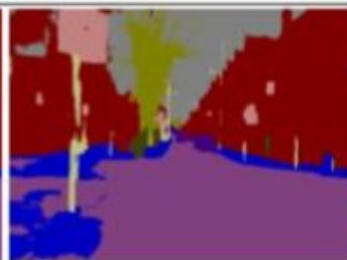
CamVid Results	IoU Accuracy
DenseNet (State of the art baseline)	67.1
+ Aleatoric Uncertainty	67.4
+ Epistemic Uncertainty	67.2
+ Aleatoric & Epistemic	67.5



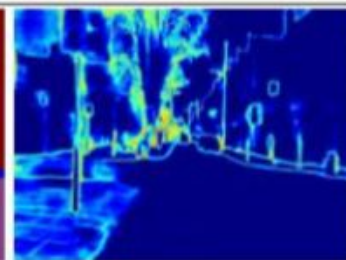
(a) Input Image



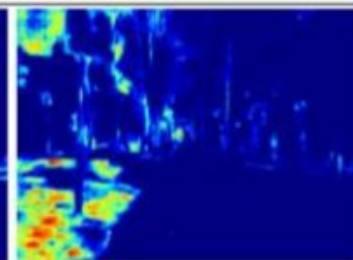
(b) Ground Truth



(c) Semantic Segmentation



(d) Aleatoric Uncertainty



(e) Epistemic Uncertainty