# Understanding models via their training data
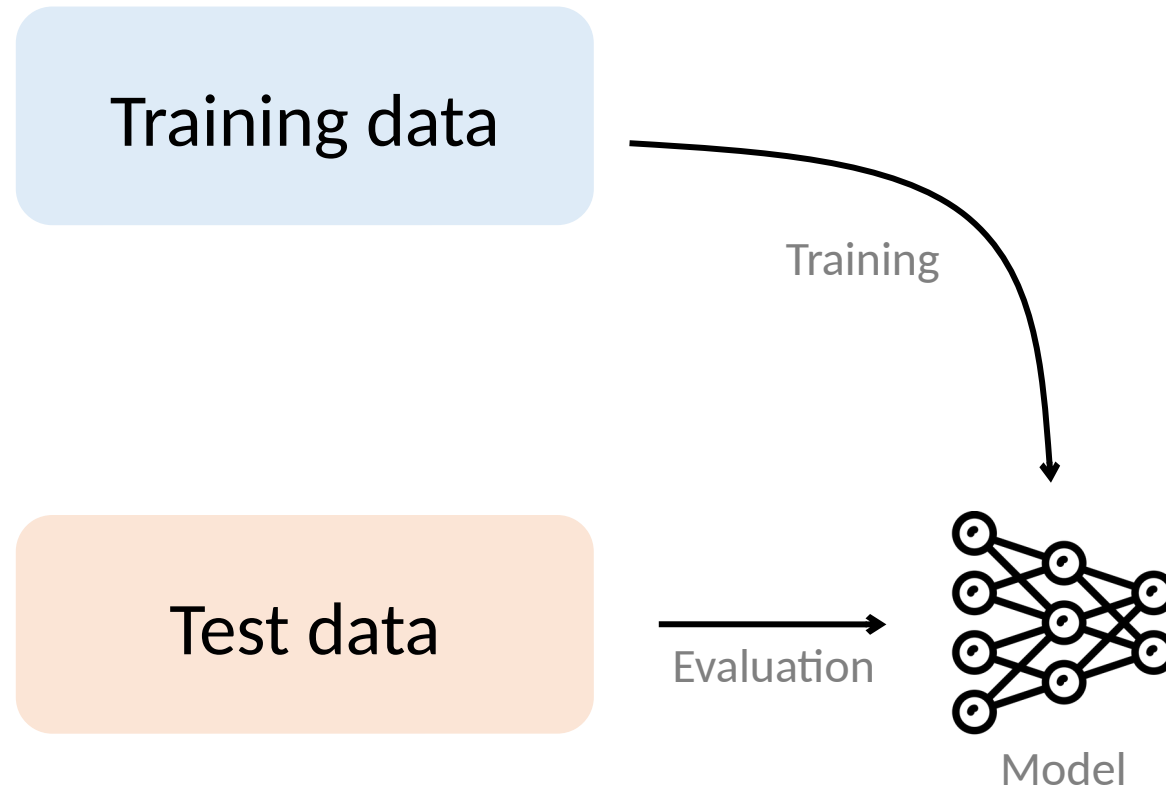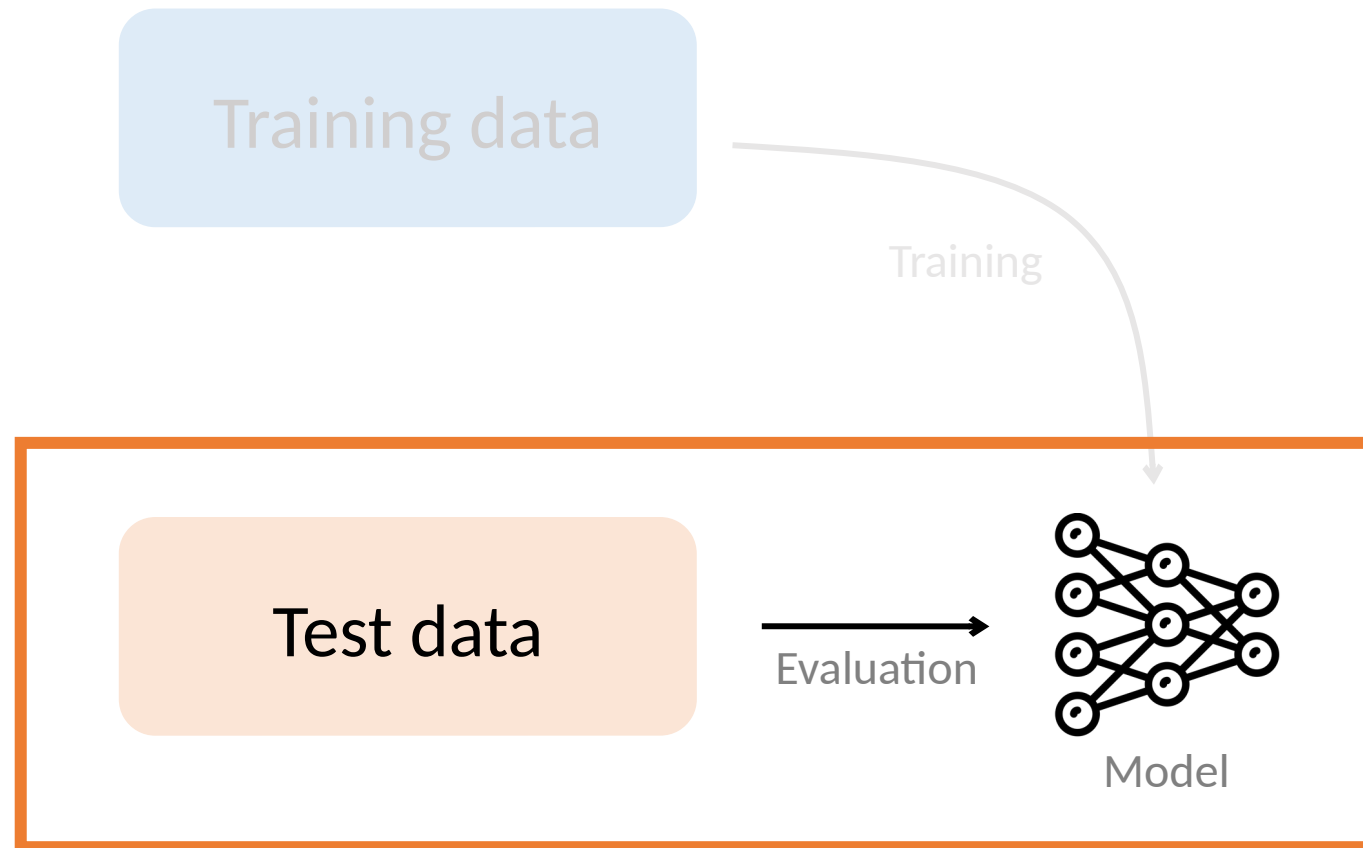
**Understanding black-box predictions via influence functions.**
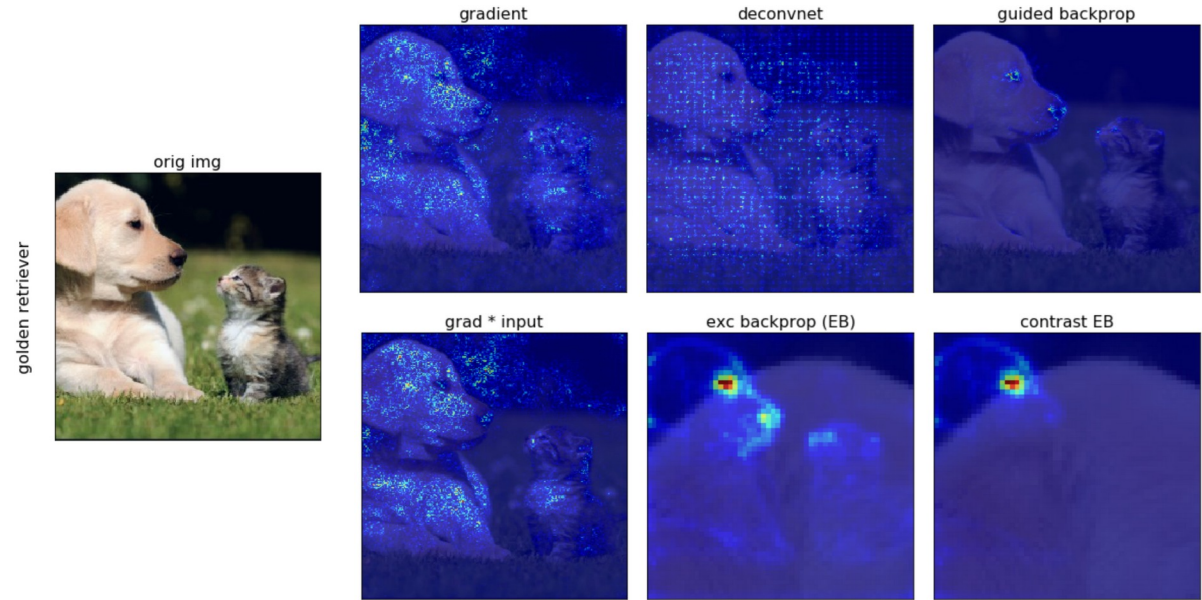Koh and Liang, ICML 2017.

Slides by Pang Wei Koh

# Prior work: Focus on test data

# Prior work: Focus on test data

# Which parts of the test input most affect the model's prediction?

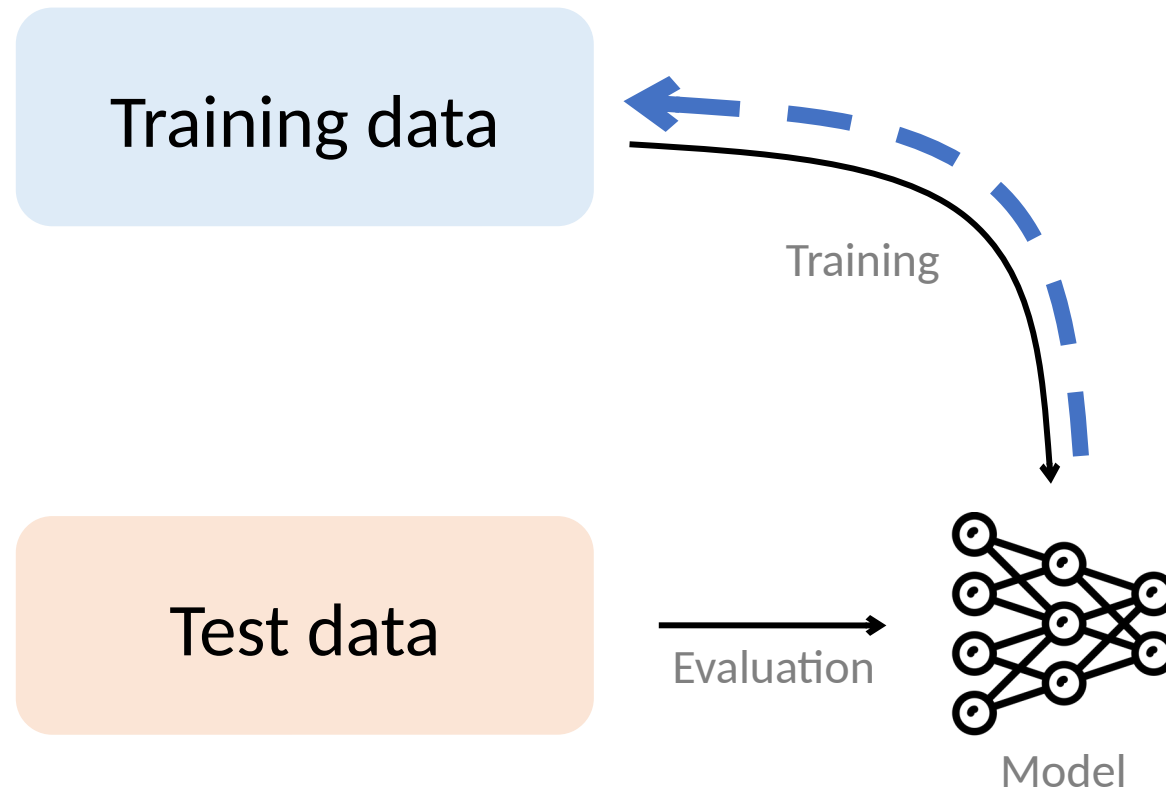

[Ruth Fong's interpretability tutorial, 2019]

Class activation maps [Zhou et al., 2016]
Concept activation vectors [Kim et al., 2018]
DeConvNet [Zeiler & Fergus, 2014]
Deep Taylor decomposition [Montavon et al., 2017]
DeepLIFT [Shrikumar et al., 2017]
Deletion game [Fong & Vedaldi, 2017]
Generalized additive models [Caruana et al., 2015]
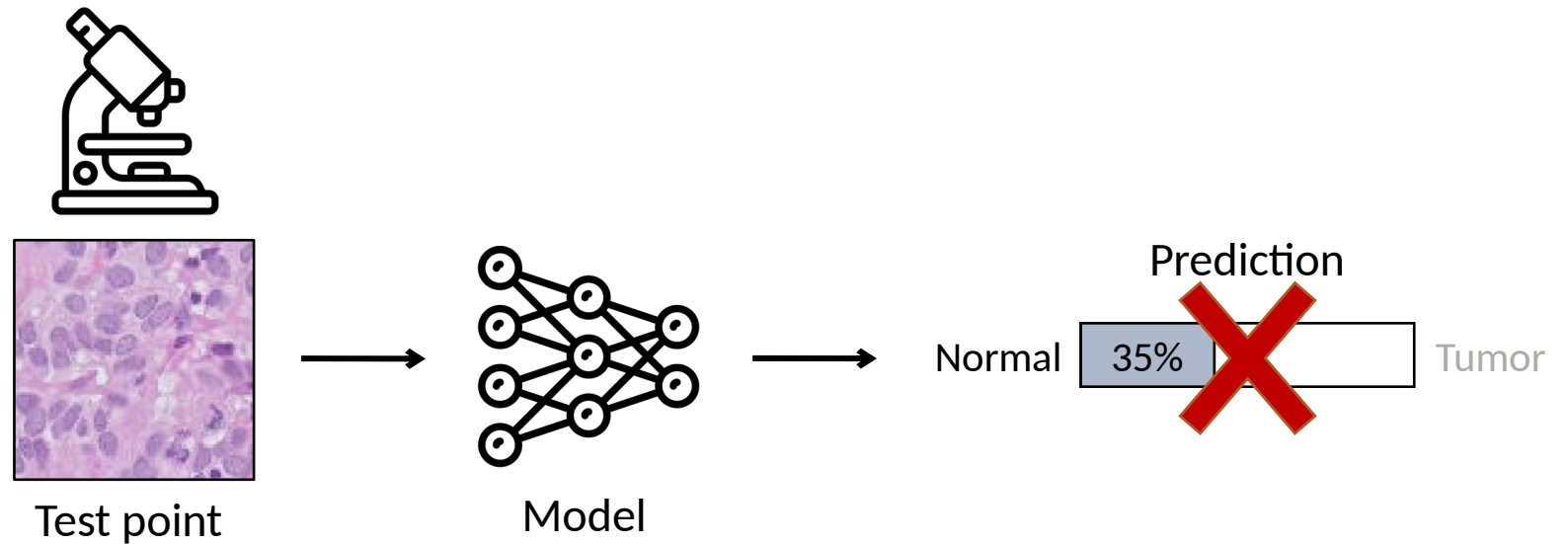Gradients [Baehrens et al., 2010]
...

Guided BackProp [Springenberg et al., 2015]
Grad-CAM [Selvaraju et al., 2016]
Integrated Gradients [Sundararajan et al., 2017]
Layer-wise relevance propagation [Bach et al., 2015]
LIME [Ribeiro et al., 2016]
Saliency maps [Simonyan et al., 2014]
SHAP [Lundberg et al., 2017]
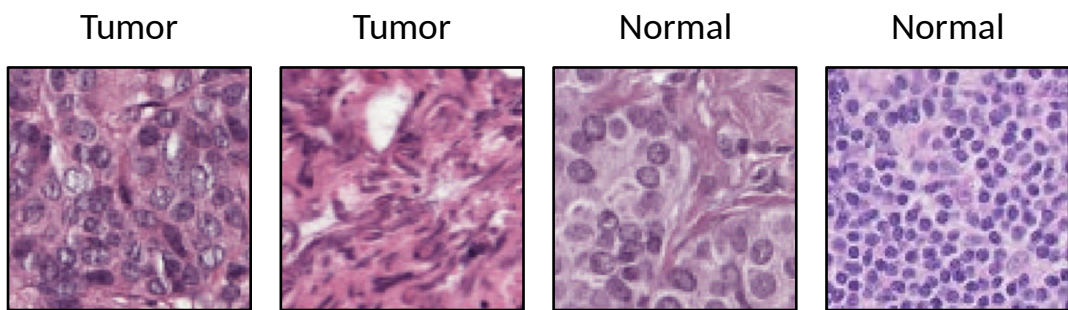SmoothGrad [Smilkov et al., 2017]
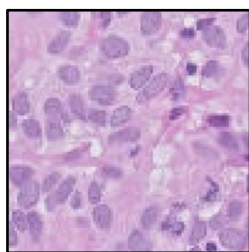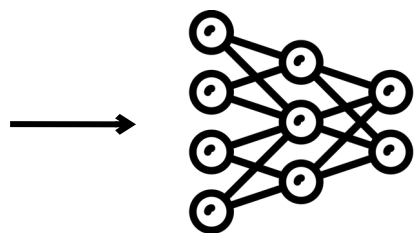...

# Our work: Link model to training data

Training data

Training

Test data

Evaluation

Model

# Example: Dataset debugging

Tumor  Tumor  Normal  Normal

Training data

Training

Test point

Model

Prediction

Normal  35%  Tumor

Tumor   Tumor   **Tumor** ~~Normal~~   Normal

Training data

Training

Test point   Model

Prediction

Normal   35%   Tumor

Tumor    Tumor    Normal    Normal

Leave-one-out approach

Training data

[Quenouille,1956; Tukey, 1958]

Training

Test point

Model

Prediction

Normal    35%    Tumor

Tumor · Tumor · Normal · Normal

Training data

Training

Test point

Model

Prediction

Normal 35% Tumor

-2%

Normal 33% Tumor

Tumor Tumor Normal Normal

Training data

Training

Test point

Model

Prediction

Normal 35% Tumor

+40%

Normal 75% Tumor
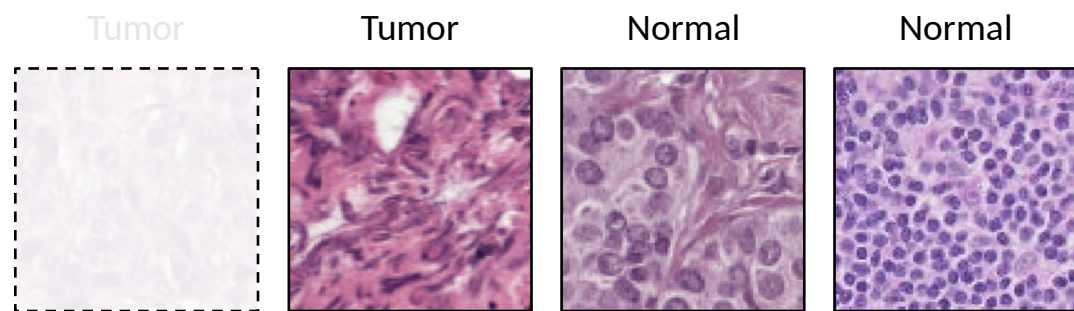
Training data

Tumor | Tumor | Normal | Normal

Weights: 1 1 1 1

Training

Test point

Model

Training data

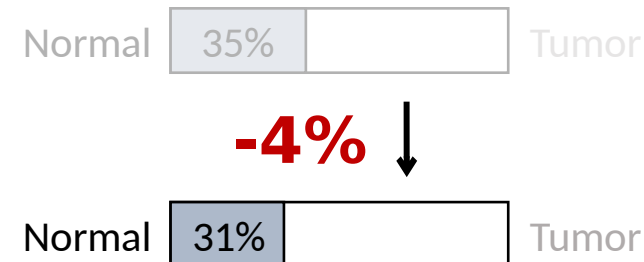Tumor　　Tumor　　Normal　　Normal

Weights　　**0**　　1　　1　　1

Training

Test point → Model

# The influence function approximation

$$z_{test}, \widehat{\theta}\Big) \approx$$

Loss on
(original)

# The influence function approximation

Change in loss on  after
removing

$$\ell\left(z_{test}, \widehat{\theta}_{-z_{train}}\right) - \ell\left(z_{test}, \widehat{\theta}\right) \approx$$

Loss on
(after removing )

Loss on
(original)

# The influence function approximation

Change in loss on  after
removing

$$\ell\left(z_{test}, \widehat{\theta}_{-z_{train}}\right) - \ell\left(z_{test}, \widehat{\theta}\right) \approx$$

# The influence function approximation

Change in loss on after removing

$$\ell\left(z_{test},\widehat{\theta}_{-z_{train}}\right) - \ell\left(z_{test},\widehat{\theta}\right) \approx \nabla_{\theta}\ell\left(z_{test},\widehat{\theta}\right)^{T} H_{\widehat{\theta}}^{-1} \nabla_{\theta}\ell(z_{train},\widehat{\theta})$$

Gradient of loss on

Gradient of loss on

Inverse of the Hessian

# The influence function approximation

Doesn't require retraining!

Change in loss on after removing

$$\ell\left(z_{test}, \hat{\theta}_{-z_{train}}\right) - \ell\left(z_{test}, \hat{\theta}\right) \approx \nabla_\theta \ell\left(z_{test}, \hat{\theta}\right)^T H_{\hat{\theta}}^{-1} \nabla_\theta \ell(z_{train}, \hat{\theta})$$

Gradient of loss on

Inverse of the Hessian

Gradient of loss on

# The influence function approximation

Effect of other training points

Model's representation of

Model's representation of

Change in loss on after removing

$$\ell\left(z_{test}, \widehat{\theta}_{-z_{train}}\right) - \ell\left(z_{test}, \widehat{\theta}\right) \approx \nabla_\theta \ell\left(z_{test}, \widehat{\theta}\right)^T H_{\widehat{\theta}}^{-1} \nabla_\theta \ell(z_{train}, \widehat{\theta})$$

Gradient of loss on

Gradient of loss on

Inverse of the Hessian

# A little bit more technical details ...

Training data $\longrightarrow \widehat{\theta} \longrightarrow \ell\left(z_{test}, \widehat{\theta}\right)$

Training data

| Tumor | Tumor | Normal |
|-------|-------|--------|

Weights $\quad +\epsilon \quad\quad 1 \quad\quad 1$

$z$

$$\hat{\theta}_{\epsilon,z} \stackrel{\text{def}}{=} \arg\min_{\theta\in\Theta} \frac{1}{n}\sum_{i=1}^{n} L(z_i, \theta) + \epsilon L(z, \theta)$$

# A little bit more technical details ...

Training data $\longrightarrow \widehat{\mathcal{O}} \longrightarrow \ell\left(z_{test}, \widehat{\theta}\right)$

$$\frac{d\hat{\theta}_{\epsilon,z}}{d\epsilon}\bigg|_{\epsilon=0} \qquad \frac{dL(z_{\text{test}}, \hat{\theta}_{\epsilon,z})}{d\epsilon}\bigg|_{\epsilon=0}$$

$$= \nabla_{\theta} L(z_{\text{test}}, \hat{\theta})^{\top} \frac{d\hat{\theta}_{\epsilon,z}}{d\epsilon}\bigg|_{\epsilon=0}$$

$$\hat{\theta}_{\epsilon,z} \stackrel{\text{def}}{=} \arg\min_{\theta\in\Theta} \frac{1}{n}\sum_{i=1}^{n} L(z_i, \theta) + \epsilon L(z, \theta)$$

# From classical to modern settings

Jaeckel, 1972. The infinitesimal jackknife.

Hampel, 1974. The influence curve and its role in robust estimation.

Cook, 1977. Detection of influential observations in linear regression.

…

# From classical to modern settings

Small datasets $\longrightarrow$ Large datasets
Low-dimensional         High-dimensional

Difficult to compute
$\downarrow$

$$\nabla_\theta \ell\left(z_{test}, \widehat{\theta}\right)^T \boldsymbol{H}_{\widehat{\boldsymbol{\theta}}}^{-1} \nabla_\theta \ell\left(z_{train}, \widehat{\theta}\right)$$

We use tools from
2nd-order optimization & stochastic estimation

[Pearlmutter, 1994; Martens, 2010, Agarwal et al., 2017]

# Removing single points

Logistic regression (MNIST)



**Estimated change in test loss (via influence functions)**

**Actual change in test loss**

Each point represents removing one training example

# Understanding Models via Their Training Data



**Where can you apply influence functions?**

# Impact

| **Robustness** | **Applications** | **Fairness & security** |
|:---:|:---:|:---:|
| Training set biases | Data distillation | |
| [Ren et al., 2018] | [Wang et al., 2020] | Algorithmic bias |
| | | [Verma et al., 2021] |
| Inference reliability | Data valuation | |
| [Broderick et al., 2021] | [Jia et al., 2019] | Data labor |
| | | [Arrieta-Ibarra, 2018] |
| Cross-validation | Active learning | |
| [Stephenson et al., 2020] | [Gudovskiy et al., 2020] | Privacy |
| | | [Shokri et al., 2021] |
| Memorization | Data debugging | |
| [Feldman, 2019] | [Guo et al., 2021] | Data poisoning |
| | | [Chen et al., 2017] |

# Limitations

- What assumptions on the models are needed to calculate influence functions?
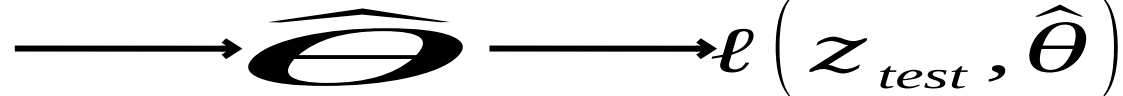
- Can you calculate influence functions for a neural network with the recipe described earlier?

$$\left. \frac{dL(z_{\text{test}}, \hat{\theta}_{\epsilon,z})}{d\epsilon} \right|_{\epsilon=0}$$

Training data $\longrightarrow \widehat{\theta} \longrightarrow \ell\left(z_{test}, \widehat{\theta}\right)$

$$\hat{\theta}_{\epsilon,z} \overset{\text{def}}{=} \arg\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} L(z_i, \theta) + \epsilon L(z, \theta)$$

# Limitations

- What assumptions on the models are needed to calculate influence functions?

- Can you calculate influence functions for a neural network with the recipe described earlier?

If Influence Functions are the Answer, Then What is the Question?

Juhan Bae*[†], Nathan Ng[†‡], Alston Lo[†], Marzyeh Ghassemi[‡], Roger Grosse[†]

Non-convexity     Non-convergence

$$\hat{\theta}_{\epsilon,z} \stackrel{\text{def}}{=} \arg\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} L(z_i, \theta) + \epsilon L(z, \theta)$$

# Summary

- Link model behavior to training data

- Efficient to calculate

- Many applications


- Limitations when applying to complex models



Training data

Test data

Influence functions

Training

Evaluation