

# The Implicit and Explicit Regularization Effects of Dropout

Bayesian Inference and Implicit Regularization

Colin Wei <sup>1</sup>   Sham Kakade <sup>2,3</sup>   Tengyu Ma <sup>1</sup>

Presenters: Kelvin Wong, Siva Manivasagam and Amanjit Singh Kainth

<sup>1</sup>Stanford University   <sup>2</sup>Microsoft Research

<sup>3</sup>University of Washington

# Dropout<sup>1</sup>

- Regularizes deep networks, especially in vision and language tasks
- Sets a random subset of the activations in each network to 0:

$$\eta = \begin{cases} -1 & \text{with probability } q \\ \frac{q}{1-q} & \text{with probability } 1 - q \end{cases}$$

$$h_{\text{drop}} = (\vec{1} + \eta) \odot h$$

---

<sup>1</sup>Srivastava et al. (2014)

# Disentangling Explicit and Implicit Regularization

## Empirical Study

- Why does it work well? Wei et al. (2020) analyze two effects of Dropout.
- Explicit regularization: an explicit change in performance by modifying the objective function ( $\mathcal{J}(x)$  vs  $\mathcal{J}_{\text{drop}}(x)$ )
- Implicit regularization: an implicit change in performance due to stochastic optimization (similar to SGD batch size)
- Let's take a look at how we can disentangle these Dropout effects in training → [Google Colab](#)

# Disentangling Explicit and Implicit Regularization

## Notation

- Standard objective  $\mathcal{J}(x) =: \mathcal{L} \circ F(x)$ , decompose  $F(x) = F_i(h_i(x))$

# Disentangling Explicit and Implicit Regularization

## Notation

- Standard objective  $\mathcal{J}(x) =: \mathcal{L} \circ F(x)$ , decompose  $F(x) = F_i(h_i(x))$
- $i$ -th partial objective  $\mathcal{J}_i(h) =: \mathcal{L} \circ F_i(h)$  (note  $\mathcal{J}(x) = \mathcal{J}_i(h_i(x))$ )

# Disentangling Explicit and Implicit Regularization

## Notation

- Standard objective  $\mathcal{J}(x) =: \mathcal{L} \circ F(x)$ , decompose  $F(x) = F_i(h_i(x))$
- $i$ -th partial objective  $\mathcal{J}_i(h) =: \mathcal{L} \circ F_i(h)$  (note  $\mathcal{J}(x) = \mathcal{J}_i(h_i(x))$ )
- Dropout objective (assuming on the  $i$ -th hidden layer)

$$\begin{aligned}\mathcal{J}_{\text{drop}}(x, \eta) &=: \mathcal{L} \circ F(x, \eta) = \mathcal{L} \circ F_i(h_i(x) + \delta) \\ &= \mathcal{J}_i(h_i(x) + \delta), \quad \delta \triangleq \eta \odot h_i(x)\end{aligned}$$

$$\mathcal{J}_{\text{drop}}(x) =: \mathbb{E}_{\eta}[\mathcal{J}_{\text{drop}}(x, \eta)] \quad (\text{Expected Dropout objective})$$

# Disentangling Explicit and Implicit Regularization

## Notation

- Standard objective  $\mathcal{J}(x) =: \mathcal{L} \circ F(x)$ , decompose  $F(x) = F_i(h_i(x))$
- $i$ -th partial objective  $\mathcal{J}_i(h) =: \mathcal{L} \circ F_i(h)$  (note  $\mathcal{J}(x) = \mathcal{J}_i(h_i(x))$ )
- Dropout objective (assuming on the  $i$ -th hidden layer)

$$\begin{aligned}\mathcal{J}_{\text{drop}}(x, \eta) &=: \mathcal{L} \circ F(x, \eta) = \mathcal{L} \circ F_i(h_i(x) + \delta) \\ &= \mathcal{J}_i(h_i(x) + \delta), \quad \delta \triangleq \eta \odot h_i(x)\end{aligned}$$

$$\mathcal{J}_{\text{drop}}(x) =: \mathbb{E}_{\eta}[\mathcal{J}_{\text{drop}}(x, \eta)] \quad (\text{Expected Dropout objective})$$

- **Explicit Regularization:** Discrepancy between (expected) training (with Dropout) and standard objectives
- **Implicit Regularization:** Stochasticity-induced, measured by the fluctuation of stochastic gradient dropout around its mean

# Explicit Regularization

## Taylor Expansion and Gauss-Newton Approximation

- Second order approximation to  $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$

$$\mathcal{J}_i(h_i(x) + \delta) \approx \mathcal{J}(x) + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle + \frac{1}{2} \delta^\top \nabla^2 \mathcal{J}_i(h_i(x)) \delta$$

$$\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x) \approx \frac{1}{2} \langle \nabla^2 \mathcal{J}_i(h_i(x)), \mathbb{E}_\eta[\delta \delta^\top] \rangle$$

# Explicit Regularization

## Taylor Expansion and Gauss-Newton Approximation

- Second order approximation to  $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$

$$\mathcal{J}_i(h_i(x) + \delta) \approx \mathcal{J}(x) + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle + \frac{1}{2} \delta^\top \nabla^2 \mathcal{J}_i(h_i(x)) \delta$$

$$\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x) \approx \frac{1}{2} \langle \nabla^2 \mathcal{J}_i(h_i(x)), \mathbb{E}_\eta[\delta \delta^\top] \rangle$$

- Recall  $\delta \triangleq \eta \odot h_i(x)$ , can show that  $\mathbb{E}_\eta[\delta \delta^\top] \propto \text{diag}(h_i(x)^{\odot 2})$

# Explicit Regularization

## Taylor Expansion and Gauss-Newton Approximation

- Second order approximation to  $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$

$$\mathcal{J}_i(h_i(x) + \delta) \approx \mathcal{J}(x) + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle + \frac{1}{2} \delta^\top \nabla^2 \mathcal{J}_i(h_i(x)) \delta$$

$$\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x) \approx \frac{1}{2} \langle \nabla^2 \mathcal{J}_i(h_i(x)), \mathbb{E}_\eta[\delta \delta^\top] \rangle$$

- Recall  $\delta \triangleq \eta \odot h_i(x)$ , can show that  $\mathbb{E}_\eta[\delta \delta^\top] \propto \text{diag}(h_i(x)^{\odot 2})$
- Use Gauss-Newton approximation on hessian

$$\nabla^2 \mathcal{J}_i(h_i(x)) \approx \nabla F_i(h_i(x))^\top \overbrace{\nabla^2 \mathcal{L}(F(x))}^{H_{\text{out}}(x)} \overbrace{\nabla F_i(h_i(x))}^{J_{F_i}(x)}$$

# Explicit Regularization

## Taylor Expansion and Gauss-Newton Approximation

- Second order approximation to  $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$

$$\mathcal{J}_i(h_i(x) + \delta) \approx \mathcal{J}(x) + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle + \frac{1}{2} \delta^\top \nabla^2 \mathcal{J}_i(h_i(x)) \delta$$

$$\mathbb{E}_\eta [\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x) \approx \frac{1}{2} \langle \nabla^2 \mathcal{J}_i(h_i(x)), \mathbb{E}_\eta [\delta \delta^\top] \rangle$$

- Recall  $\delta \triangleq \eta \odot h_i(x)$ , can show that  $\mathbb{E}_\eta [\delta \delta^\top] \propto \text{diag}(h_i(x)^{\odot 2})$
- Use Gauss-Newton approximation on hessian

$$\nabla^2 \mathcal{J}_i(h_i(x)) \approx \nabla F_i(h_i(x))^\top \overbrace{\nabla^2 \mathcal{L}(F(x))}^{H_{\text{out}}(x)} \overbrace{\nabla F_i(h_i(x))}^{J_{F_i}(x)}$$

- Final Form of Explicit Regularizer (multi-layer)

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top H_{\text{out}}(x) J_{F_i}(x), \text{diag}(h_i(x)^{\odot 2}) \right\rangle$$

# Explicit Regularization

Cross-entropy for classification tasks

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top H_{\text{out}}(x) J_{F_i}(x), \text{diag}\left(h_i(x)^{\odot 2}\right) \right\rangle$$

- Full  $H_{\text{out}}(x) \in \text{Mat}(|V| \times |V|)$ ,  $|V|$  can be quite large!

# Explicit Regularization

## Cross-entropy for classification tasks

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top H_{\text{out}}(x) J_{F_i}(x), \text{diag}\left(h_i(x)^{\odot 2}\right) \right\rangle$$

- Full  $H_{\text{out}}(x) \in \text{Mat}(|V| \times |V|)$ ,  $|V|$  can be quite large!
- For the cross-entropy loss  $\mathcal{L}_y^{\text{ce}}(\mathbf{v}) = \mathcal{L}^{\text{ce}}(\mathbf{v}, y) = -\log \text{softmax}(\mathbf{v})_y$ ,

$$H_{\text{out}}^{\text{ce}}(x) = \mathbb{E}_{\hat{y} \sim \text{softmax}(F(x))} \left[ \nabla \mathcal{L}_{\hat{y}}^{\text{ce}}(F(x))^\top \nabla \mathcal{L}_{\hat{y}}^{\text{ce}}(F(x)) \right]$$

# Explicit Regularization

## Cross-entropy for classification tasks

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top H_{\text{out}}(x) J_{F_i}(x), \text{diag}\left(h_i(x)^{\odot 2}\right) \right\rangle$$

- Full  $H_{\text{out}}(x) \in \text{Mat}(|V| \times |V|)$ ,  $|V|$  can be quite large!
- For the cross-entropy loss  $\mathcal{L}_y^{\text{ce}}(\mathbf{v}) = \mathcal{L}^{\text{ce}}(\mathbf{v}, y) = -\log \text{softmax}(\mathbf{v})_y$ ,

$$H_{\text{out}}^{\text{ce}}(x) = \mathbb{E}_{\hat{y} \sim \text{softmax}(F(x))} \left[ \nabla \mathcal{L}_{\hat{y}}^{\text{ce}}(F(x))^\top \nabla \mathcal{L}_{\hat{y}}^{\text{ce}}(F(x)) \right]$$

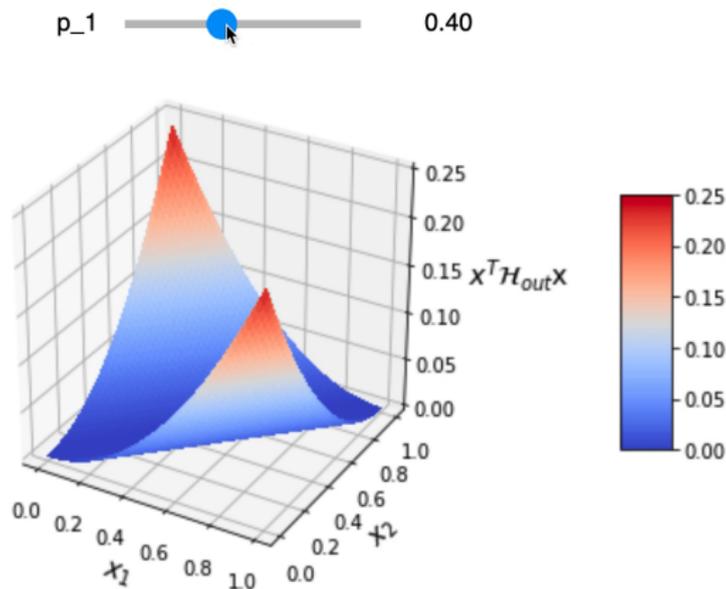
- Memory-efficient (unbiased) estimator for  $R_{\text{approx}}$

$$\hat{R}_{\text{approx}} = \left\langle \nabla \mathcal{J}_i^{\hat{y}}(h_i(x))^\top \nabla \mathcal{J}_i^{\hat{y}}(h_i(x)), \text{diag}\left(h_i(x)^{\odot 2}\right) \right\rangle$$

where  $\nabla \mathcal{J}_i^{\hat{y}}(h_i(x)) =: \nabla \mathcal{L}_{\hat{y}}^{\text{ce}}(F(x)) J_{F_i}(x)$  for  $\hat{y} \sim \text{softmax}(F(x))$

# Explicit Regularization

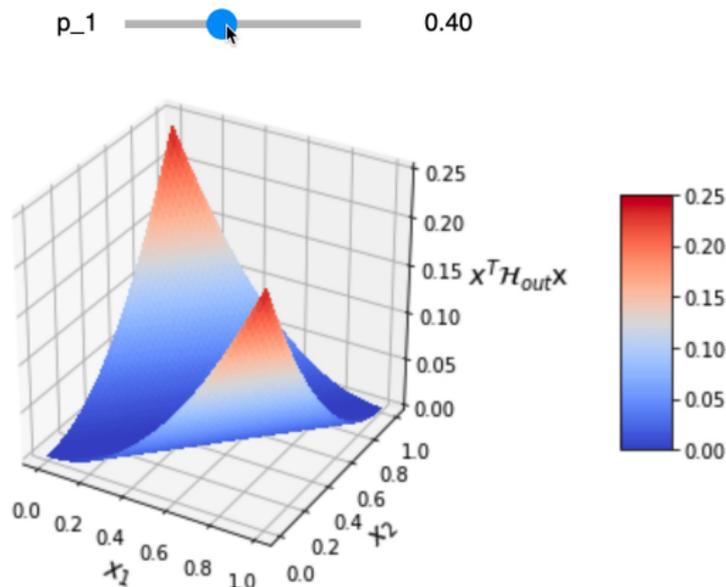
## Cross-entropy Hessian



$$H_{out}^{ce}(x) = \nabla^2 \mathcal{L}_y^{ce}(F(x)) = \text{diag}(p) - pp^T, \quad p = \text{softmax}(F(x))$$

# Explicit Regularization

## Cross-entropy Hessian



$$H_{\text{out}}^{\text{ce}}(x) = \nabla^2 \mathcal{L}_y^{\text{ce}}(F(x)) = \text{diag}(p) - pp^T, \quad p = \text{softmax}(F(x))$$

Intuition: Penalize classes with plausible but lower confidence

# Implicit Regularization

Identifying the noise in Dropout and approximating it

- Stochastic gradient noise due to Dropout:

$$\xi_{\text{drop}}(F, x, \eta) \triangleq \nabla_W \mathcal{J}_{\text{drop}}(x, \eta) - \nabla_W \mathbb{E}_{\eta'} [\mathcal{J}_{\text{drop}}(x, \eta')]$$

# Implicit Regularization

Identifying the noise in Dropout and approximating it

- Stochastic gradient noise due to Dropout:

$$\xi_{\text{drop}}(F, x, \eta) \triangleq \nabla_W \mathcal{J}_{\text{drop}}(x, \eta) - \nabla_W \mathbb{E}_{\eta'} [\mathcal{J}_{\text{drop}}(x, \eta')]$$

- Inject noise back by taking the difference in Dropout gradient samples:

$$\tilde{\xi}_{\text{drop}}(F, x, \eta_1, \eta_2) \triangleq \nabla_W [\mathcal{J}_{\text{drop}}(x, \eta_1) - \mathcal{J}_{\text{drop}}(x, \eta_2)]$$

# Implicit Regularization

## Identifying the noise in Dropout and approximating it

- Stochastic gradient noise due to Dropout:

$$\xi_{\text{drop}}(F, x, \eta) \triangleq \nabla_W \mathcal{J}_{\text{drop}}(x, \eta) - \nabla_W \mathbb{E}_{\eta'} [\mathcal{J}_{\text{drop}}(x, \eta')]$$

- Inject noise back by taking the difference in Dropout gradient samples:

$$\tilde{\xi}_{\text{drop}}(F, x, \eta_1, \eta_2) \triangleq \nabla_W [\mathcal{J}_{\text{drop}}(x, \eta_1) - \mathcal{J}_{\text{drop}}(x, \eta_2)]$$

- We saw in the Colab that we can empirically use  $\tilde{\xi}_{\text{drop}}$  to model  $\xi_{\text{drop}}$

# Implicit Regularization

Taylor Expansion (Linear approximation)

$$\tilde{\xi}_{\text{drop}} \left( F, x, \eta_i^{(1)}, \eta_i^{(2)} \right) \triangleq \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(1)} \right) - \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(2)} \right) \quad (1)$$

- $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$ , where  $\delta \triangleq \eta \odot h_i(x)$  is the perturbation on the output of the  $i$ -th hidden layer

# Implicit Regularization

## Taylor Expansion (Linear approximation)

$$\tilde{\xi}_{\text{drop}} \left( F, x, \eta_i^{(1)}, \eta_i^{(2)} \right) \triangleq \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(1)} \right) - \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(2)} \right) \quad (1)$$

- $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$ , where  $\delta \triangleq \eta \odot h_i(x)$  is the perturbation on the output of the  $i$ -th hidden layer
- First order approximation:  $\mathcal{J}_i(h_i(x) + \delta) \approx \overbrace{\mathcal{J}_i(h_i(x))}^{\mathcal{J}(x)} + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle$

# Implicit Regularization

## Taylor Expansion (Linear approximation)

$$\tilde{\xi}_{\text{drop}} \left( F, x, \eta_i^{(1)}, \eta_i^{(2)} \right) \triangleq \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(1)} \right) - \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(2)} \right) \quad (1)$$

- $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$ , where  $\delta \triangleq \eta \odot h_i(x)$  is the perturbation on the output of the  $i$ -th hidden layer
- First order approximation:  $\mathcal{J}_i(h_i(x) + \delta) \approx \overbrace{\mathcal{J}_i(h_i(x))}^{\mathcal{J}(x)} + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle$
- Thus above simplifies to:

$$\approx \nabla_W \left\langle \nabla \mathcal{J}_i(h_i(x)), (\eta_i^{(1)} - \eta_i^{(2)}) \odot h_i(x) \right\rangle \quad (2)$$

# Implicit Regularization

## Taylor Expansion (Linear approximation)

$$\tilde{\xi}_{\text{drop}} \left( F, x, \eta_i^{(1)}, \eta_i^{(2)} \right) \triangleq \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(1)} \right) - \nabla_W \mathcal{J}_{\text{drop}} \left( x, \eta_i^{(2)} \right) \quad (1)$$

- $\mathcal{J}_{\text{drop}}(x, \eta) = \mathcal{J}_i(h_i(x) + \delta)$ , where  $\delta \triangleq \eta \odot h_i(x)$  is the perturbation on the output of the  $i$ -th hidden layer
- First order approximation:  $\mathcal{J}_i(h_i(x) + \delta) \approx \overbrace{\mathcal{J}_i(h_i(x))}^{\mathcal{J}(x)} + \langle \nabla \mathcal{J}_i(h_i(x)), \delta \rangle$
- Thus above simplifies to:

$$\approx \nabla_W \left\langle \nabla \mathcal{J}_i(h_i(x)), (\eta_i^{(1)} - \eta_i^{(2)}) \odot h_i(x) \right\rangle \quad (2)$$

- Replace  $(\eta_i^{(1)} - \eta_i^{(2)})$  with  $\sqrt{2}\eta_i$  to maintain same noise covariance

# Implicit Regularization

## Intuition

$$\xi_{\text{approx}}(F, x, \{\eta_i\}) \triangleq \nabla_W \left( \sum_i \langle \nabla \mathcal{J}_i(h_i(x)), (\eta_i \odot h_i(x)) \rangle \right)$$

- Intuition: provides data-dependent stability during training

# Explicit vs Implicit Regularization

## Summary

- **Explicit** (uses *loss hessian* and *model jacobian*)

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top H_{\text{out}}(x) J_{F_i}(x), \text{diag}\left(h_i(x)^{\odot 2}\right) \right\rangle$$

- **Implicit** (uses *loss jacobian* and *model jacobian*)

$$\xi_{\text{approx}}(F, x, \{\eta_i\}) \triangleq \nabla_W \left( \sum_i \langle \nabla \mathcal{J}_i(h_i(x)), (\eta_i \odot h_i(x)) \rangle \right)$$

- For the cross-entropy loss  $\mathcal{L}^{\text{ce}}$ , both penalize the loss and model jacobians  $\nabla \mathcal{J}_i(h_i(x)) = \nabla \mathcal{L}_y^{\text{ce}}(F(x)) J_{F_i}(x)$  (in practice)

# Colab

Let's take a look at these regularizers in action

# Conclusion

## Some key takeaways

- Dropout induces distinct regularization effects:

# Conclusion

## Some key takeaways

- Dropout induces distinct regularization effects:
  - **Explicit Regularization:** Induced by change in training objective.

$$\mathbb{E}_x[\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x)] \quad (3)$$

# Conclusion

## Some key takeaways

- Dropout induces distinct regularization effects:
  - **Explicit Regularization:** Induced by change in training objective.

$$\mathbb{E}_x[\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x)] \quad (3)$$

- **Implicit Regularization:** Induced by stochastic approximation of Dropout training objective.

$$\nabla_W \mathcal{J}_{\text{drop}}(x, \eta) - \nabla_W \mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] \quad (4)$$

# Conclusion

## Some key takeaways

- Dropout induces distinct regularization effects:
  - **Explicit Regularization:** Induced by change in training objective.

$$\mathbb{E}_x[\mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] - \mathcal{J}(x)] \quad (3)$$

- **Implicit Regularization:** Induced by stochastic approximation of Dropout training objective.

$$\nabla_W \mathcal{J}_{\text{drop}}(x, \eta) - \nabla_W \mathbb{E}_\eta[\mathcal{J}_{\text{drop}}(x, \eta)] \quad (4)$$

- We can derive analytical forms of these regularizers, revealing:
  - Explicit regularization encourages output Jacobians and hidden layers to be small according to output Hessian.
  - Implicit regularization penalizes loss Jacobian and hidden layers.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.
  - It is unclear why gradient noise helps generalization.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.
  - It is unclear why gradient noise helps generalization.
  - On larger datasets and/or other model architectures, the benefits of Dropout's implicit regularization effect is missing.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.
  - It is unclear why gradient noise helps generalization.
  - On larger datasets and/or other model architectures, the benefits of Dropout's implicit regularization effect is missing.
  - In fact, implicit regularization sometimes makes things worse.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.
  - It is unclear why gradient noise helps generalization.
  - On larger datasets and/or other model architectures, the benefits of Dropout's implicit regularization effect is missing.
  - In fact, implicit regularization sometimes makes things worse.
- The derived analytical regularizers are totally impractical.

# Conclusion

## Limitations and outlook

- Interpretation of Dropout's implicit regularizer remains opaque.
  - It is unclear why gradient noise helps generalization.
  - On larger datasets and/or other model architectures, the benefits of Dropout's implicit regularization effect is missing.
  - In fact, implicit regularization sometimes makes things worse.
- The derived analytical regularizers are totally impractical.
- But this work provides intuition on Dropout's inner workings, hopefully leading to better, more principled regularizers.

# Conclusion

Thanks for listening! Questions?

## References I

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15 (56):1929–1958, 2014. URL <http://jmlr.org/papers/v15/srivastava14a.html>.

Colin Wei, Sham Kakade, and Tengyu Ma. The implicit and explicit regularization effects of dropout. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 10181–10192. PMLR, 13–18 Jul 2020. URL <http://proceedings.mlr.press/v119/wei20d.html>.

# Appendix: Explicit Regularization

Additional notes: Using the loss hessian is important

$$R_{\text{approx}}(F, x) \triangleq \sum_i \left\langle J_{F_i}(x)^\top \overbrace{H_{\text{out}}(x)}^{\text{substitute}} J_{F_i}(x), \text{diag}(h_i(x)^{\odot 2}) \right\rangle$$

Training Method	Best Val. Ppl.
$\ell_2$ reg (tuned)	112.04
<b>I</b>	108.76
$\tilde{R}_{\text{approx}}$ (tuned)	84.06
$R_{\text{approx}}$	84.52

**Table:** Effect of explicit regularizer only

- Using identity **I** (classical GN matrix) removes most of Dropout benefits
- Using loss Jacobian:  $\tilde{R}_{\text{approx}}(F, x) \triangleq \sum_i \nabla \mathcal{J}_i(h_i(x)) \text{diag}(h_i(x)^{\odot 2}) \nabla \mathcal{J}_i(h_i(x))^\top$   
 $\rightarrow$  differs from  $R_{\text{approx}}$  implementation by using label  $y$  instead of sample  $\hat{y}$

# Appendix: Implicit Regularization

Additional notes: Dataset dependent

For large datasets, k-Dropout performs the same or better than the original Dropout.

**Table:** Experimental results on the full WikiText-103 dataset for QRNN architecture.

Training Method	Best Val. Ppl.
DROPOUT <sub>1</sub>	34.24
DROPOUT <sub>2</sub>	33.35
DROPOUT <sub>4</sub>	32.74
DROPOUT <sub>8</sub>	32.78

More study required to better understand implicit regularization - different effects might be inter-connected.