Introduction to Machine Learning

Regularization Intuition for L2 Regularization in Non-Linear Models

X X X

Learning goals

Understand how regularization and parameter shrinkage can be beneficial to non-linear models

SUMMARY: REGULARIZED RISK MINIMIZATION

If we should define (supervised) ML in only one line, this might be it:

$$
\min_{\theta} \mathcal{R}_{reg}(\theta) = \min_{\theta} \left(\sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} | \theta\right)\right) + \lambda \cdot J(\theta) \right)
$$

We can choose for a task at hand:

- **•** the **hypothesis space** of *f*, which determines how features can influence the predicted *y*
- \bullet the loss function *L*, which measures how errors should be treated
- **•** the **regularization** $J(\theta)$, which encodes our inductive bias and preference for certain simpler models

By varying these choices one can construct a huge number of different ML models. Many ML models follow this construction principle or can be interpreted through the lens of regularized risk minimization.

REGULARIZATION IN NEURAL NETWORKS

For neural networks, the regularized loss function is:

$$
\mathcal{R}_{reg}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} | \boldsymbol{\theta}\right)\right) + \lambda \cdot J(\boldsymbol{\theta})
$$

where:

- $L(f(x_i; \theta), y_i)$ is the loss function.
- $f(x_i; \theta)$ is the neural network's prediction.
- $J(\theta)$ is the regularization term (e.g., $\|\theta\|_2^2$ for L2 regularization).
- \bullet λ is the regularization parameter.

Bias: Regularization increases bias because it adds a constraint on the network parameters, preventing them from fitting the training data perfectly.

Variance: Regularization decreases variance by limiting the network parameters' magnitudes, reducing sensitivity to the training data's noise.

FORMAL BOUNDS

Consider a neural network with parameters θ trained with L2 regularization:

$$
\|\boldsymbol{\theta}\|_2^2 = \sum_{j=1}^p \theta_j^2
$$

The regularized loss function becomes:

$$
\mathcal{R}_{reg}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} | \boldsymbol{\theta}\right)\right) + \lambda \|\boldsymbol{\theta}\|_2^2
$$

To bound the variance term, note that the regularization term $\lambda\|\bm{\theta}\|_2^2$ constrains the parameters:

• Without regularization ($\lambda = 0$), the parameters can grow large, leading to high variance.

FORMAL BOUNDS / 2

• With regularization ($\lambda > 0$), the parameters are constrained, reducing variance.

Formally, the variance of the model can be bounded as follows:

$$
\textsf{Var}(\hat{\theta}_{\textsf{Reg}}) \leq \frac{\sigma^2}{\lambda}
$$

where σ^2 is the noise variance. As λ increases, the bound on the variance decreases.

DERIVING THE BOUND FOR VARIANCE OF NEURAL NETWORK PREDICTIONS

To derive the bound for the variance of the parameter estimates in a neural network with L2 regularization, we follow these steps: **Neural Network with L2 Regularization:** The regularized loss function is:

$$
\mathcal{R}_{reg}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} | \boldsymbol{\theta}\right)\right) + \lambda \|\boldsymbol{\theta}\|_2^2
$$

Bias-Variance Decomposition: The mean squared error (MSE) decomposition is:

$$
E[(\hat{y} - y)^2] = Bias^2(\hat{y}) + Var(\hat{y}) + \sigma^2
$$

Step-by-Step Derivation:

Model the Neural Network Parameters: $\hat{\theta} = \theta^* + \epsilon$

DERIVING THE BOUND FOR VARIANCE OF NEURAL NETWORK PREDICTIONS / 2

Apply Regularization:

 $\hat{\theta}_{\text{Reg}} = \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^{n} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} \mid \theta\right)\right) + \lambda \|\theta\|_2^2 \right\}$

Analyzing the Variance: Var $(\hat{\theta}_{\mathsf{Reg}}) \approx (\mathit{I}(\theta) + 2 \lambda \mathit{I})^{-1} \sigma^2$

Bounding the Variance: Given the properties of the Hessian matrix *H*:

$$
\text{Var}(\hat{\theta}_{\text{Reg}}) \leq \frac{\sigma^2}{2\lambda}I
$$

The variance of the neural network prediction is bounded by:

$$
\text{Var}(f(x; \hat{\theta}_{\text{Reg}})) \leq \frac{\sigma^2}{2\lambda} \|\nabla_{\theta} f(x; \hat{\theta}_{\text{Reg}})\|^2
$$

Conclusion: Regularization reduces the variance of the parameter estimates and helps in reducing overfitting by balancing the bias and variance.

BIAS ANALYSIS IN NEURAL NETWORKS

To analyze the bias term:

Bias Term: Regularization introduces bias by shrinking the parameter estimates towards zero:

$$
Bias(f(x)) = E[f(x; \hat{\theta}_{Reg})] - f^{*}(x)
$$

Using a linear approximation:

$$
E[f(x; \hat{\theta}_{\text{Reg}})] \approx f(x; \theta^*) - \lambda \nabla_{\theta} f(x; \theta^*)^T H^{-1} \theta^*
$$

Thus, the bias is:

Bias
$$
(f(x)) = -\lambda \nabla_{\theta} f(x; \theta^*)^T H^{-1} \theta^*
$$

Combined Bias and Variance Analysis:

- **Bias:** Bias² $(f(x)) = (\lambda \nabla_{\theta} f(x; \theta^*)^T H^{-1} \theta^*)^2$
- **Variance:** $\mathsf{Var}(f(x; \hat{\theta}_{\mathsf{Reg}})) \leq \frac{\sigma^2}{2\lambda}$ $\frac{\sigma^2}{2\lambda} \| \nabla_{\boldsymbol{\theta}} f (\boldsymbol{x} ; \hat{\boldsymbol{\theta}}_{\mathsf{Reg}}) \|^2$

REDUCTION IN VARIANCE VS. INCREASE IN BIAS

To show that the reduction in variance is usually more than the increase in bias, consider:

Bias-Variance Trade-off: The MSE is decomposed as:

 $\mathsf{MSE} = \mathsf{Bias}^2(f(x)) + \mathsf{Var}(f(x)) + \sigma^2$

Change in Bias and Variance:

- **Change in Bias:** \triangle Bias² $\propto \lambda^2$
- **Change in Variance:** Δ Var $\propto -\frac{1}{\lambda}$

For small λ , the reduction in variance is significant, while the increase in bias is relatively small. The reduction in variance usually outweighs the increase in bias, leading to an overall decrease in MSE. **Conclusion:** Regularization helps in reducing the overall prediction error by balancing the bias and variance effectively.

 $\overline{\mathbf{X}}$

CRITIQUE: BIAS-VARIANCE TRADEOFF AND OPTIMIZATION

For linear models, it's well-established that some $\lambda > 0$ can balance the increase in bias against the reduction in variance, leading to a net decrease in MSE. For non-linear models, the situation is more complex:

- The relationship between model parameters θ , the regularization term, and the model output $f(x; \theta)$ is non-linear.
- The effects of changing λ on the bias and variance terms are not straightforward and depend heavily on the specific form of the non-linear model and the data distribution.

Proving analytically that there exists a $\lambda > 0$ such that the regularized model always outperforms the unregularized model in terms of MSE for general non-linear models involves:

• Detailed understanding of how changes in λ affect the bias and variance for the specific type of non-linear model.

CRITIQUE: BIAS-VARIANCE TRADEOFF AND OPTIMIZATION / 2

Possibly making assumptions about the smoothness, continuity, or differentiability of the model function *f* with respect to both *x* and θ.

CRITIQUE: CONCLUSION

In summary, while it is conceptually feasible to argue that an appropriate $\lambda > 0$ might improve the MSE by balancing bias and variance, providing a universal, formal proof for all non-linear models would require either restrictive assumptions about the models and data or a very specific setup where the non-linearities are well understood and mathematically tractable.

For practical purposes, empirical validation through techniques such as cross-validation remains a critical method to determine the optimal λ for specific non-linear models and datasets.

l X

COUNTEREXAMPLE

Chris: I think ChatGPT produced a lot of "almost correct" stuff that culminated in a globally useless derivation. A general proof for DNNs imo can not work by giving a simple counterexample.

- A diagonal linear network with one hidden layer and one output unit can be written as $f(\pmb{\mathsf{x}}|\pmb{\mathsf{u}},\pmb{\mathsf{v}}) = (\pmb{\mathsf{u}}\odot \pmb{\mathsf{v}})^\top \pmb{\mathsf{x}}$
- **•** optimizing the network with L2 regularization λ and MSE loss has multiple global minima that coincide with the lasso solution for the collapsed parameter $\theta := u \odot v$ using 2λ
- Since there is no existence theorem (of a λ^* that reduces the MSE over OLS) for lasso compared to ridge regression, there can not be one for L2 regularized DNNs in general.

 $^{\prime}$ \times

COUNTEREXAMPLE / 2

- For fully-connected linear networks using *L* weight matrices $f(x|W_1,\ldots,W_1) = W_1 \cdot \ldots \cdot W_1$ *x*, adding *L*2 regularization with λ to all *W^l* produces equivalent minma to Schatten 2/*L*-norm regularization of the the collapsed linear predictor $\overline{W}X := W_1 \cdot \ldots \cdot W_1X$ with strength $L\lambda$
- I am fairly certain there is also no existence theorem for non-convex Schatten 2/*L*-norm regularization, their success depends strongly on the low-rank nature of the problem
- For MLPs beyond linear DNNs there are also some results for the "induced regularizer" in specific cases, which is often a complex or non-analytical expression. For these, there are also no existence theorems

COUNTEREXAMPLE / 3

[Neyshabur et al., 2015](https://arxiv.org/pdf/1412.6614) derive equivalent optimization problems for *L*2 regularized shallow relu-networks:

$$
\underset{\boldsymbol{v}\in\mathbb{R}^{H},(\boldsymbol{u}_{h})_{h=1}^{H}}{\operatorname{argmin}}\left(\sum_{t=1}^{n}L\left(y_{t},\sum_{h=1}^{H}v_{h}\left[\langle\boldsymbol{u}_{h},\boldsymbol{x}_{t}\rangle\right]_{+}\right)+\frac{\lambda}{2}\sum_{h=1}^{H}\left(\left\Vert\boldsymbol{u}_{h}\right\Vert^{2}+\left|v_{h}\right|^{2}\right)\right),
$$

is the same as

$$
\underset{\mathbf{v}\in\mathbb{R}^{H},(\mathbf{u}_{h})_{h=1}^{H}}{\operatorname{argmin}}\left(\sum_{t=1}^{n}L\left(y_{t},\sum_{h=1}^{H}v_{h}\left[\langle\mathbf{u}_{h},\mathbf{x}_{t}\rangle\right]_{+}\right)+\lambda\sum_{h=1}^{H}\left|v_{h}\right|\right),
$$
\nsubject to $||\mathbf{u}_{h}||\leq 1$ $(h=1,\ldots,H).$

How can we do a general analysis of the effect of *L*2 regularization in DNNs when there are these close connections to other regularized problems for which there is no anaysis of the bias-variance trade-off and no existence theorem of an optimal $\lambda^* > 0$?