

# Generalizaiton and Regularizaiton

**Tianxiang (Adam) Gao**

September 30, 2024

# Outline

1 Statistical Learning Theory

2 Regularization

3 Hyperparameter Tune

4 Overparameterization

# Recap: Optimization in Neural Network

- MLPs are **parameterized** function  $f_{\theta}$ , where  $\theta = \{W^{\ell}, b^{\ell}\}$
- The learning problem involves solving an **optimization** problem to iteratively update the  $\theta$

$$\min_{\theta} \quad \mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(\mathbf{x}_i), \mathbf{y}_i)$$

where  $\ell$  is a **loss** function and  $\mathcal{S} := \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{\ell}$  is a **training set**.

- This optimization problem can be solved using **gradient**-based methods such as (*stochastic*) *gradient descent (SGD)*, *gradient descent with momentum*, *RMSPprop*, *Adam*, etc:

$$\theta^{+} = \theta - \eta \cdot \mathbf{v}^{+},$$

where  $\eta > 0$  is a **learning rate** and  $\mathbf{v}$  is a **search direction**.

## Recap: Hyperparameters in Neural Networks

The training process involves several key hyperparameters:

- **Loss Function**  $\ell(\cdot, \cdot)$ : Square loss, cross-entropy loss, hinge loss
- **Activation Function**  $\phi(\cdot)$ : Step, sigmoid, ReLU, tanh, GELU
- **Optimizer**: SGD, Momentum, RMSProp, Adam
- **Learning Rate** ( $\eta$ ), **Batch Size** ( $b$ ), **Epochs**
- **Network Type**: MLPs, CNNs, RNNs, Transformers, GNNs
- **Width and Depth**
- **Layers**: Normalization, pooling, dropout, softmax
- **Otherwise**: Initialization (Xavier, He),  $\ell_2$ -regularization, gradient clipping, early stop

### Key Difference: Hyperparameters vs. Trainable Parameters

- Hyperparameters are **not trainable**. Unlike weights and biases, they need to be **tuned manually**.
- Proper tuning is essential for faster convergence during training and achieving good **generalization** performance.

# Outline

## 1 Statistical Learning Theory

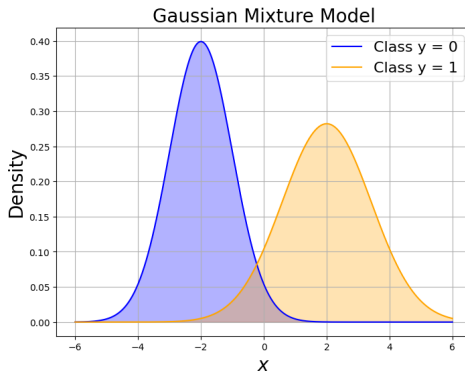
## 2 Regularization

## 3 Hyperparameter Tune

## 4 Overparameterization

# Gaussian Mixture Model

- Assume the output  $y$  follows a *discrete uniform distribution* over  $\{0, 1\}$ , meaning  $y \sim \mathcal{U}\{0, 1\}$ .
- For each value of  $y$ , the input  $x$  follows a *Gaussian distribution*:
  - When  $y = 0$ ,  $x$  follows  $x|y = 0 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ .
  - When  $y = 1$ ,  $x$  follows  $x|y = 1 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ .



- This setup defines a **(binary) Gaussian Mixture Model (GMM)**.
- Both  $x$  and  $y$  are random variables, with a joint distribution denoted as  $\mathcal{D}$ , i.e.,  $(x, y) \sim \mathcal{D}$ .

# Statistical Learning Theory (SLT)

- Assume the data  $(x, y)$  is drawn from an underlying joint distribution  $\mathcal{D}$ , i.e.,  $(x, y) \sim \mathcal{D}$ .
- The goal of learning is to find a (parameterized) function  $f$  such that:

$$f(x) \approx y$$

for "most"  $(x, y)$  pairs in a probabilistic sense.

- The **expected risk** of  $f$  is defined as:

$$R(f) := \mathbb{E}_{(x,y) \sim \mathcal{D}} [f(x) - y]^2,$$

where we use the squared loss to measure the difference between  $f(x)$  and  $y$ .

- In practice, the distribution  $\mathcal{D}$  is **unknown**.
- Instead, we collect a **random training sample**  $\mathcal{S} := \{(x_i, y_i)\}_{i=1}^n$  and compute the **empirical risk** or **training error**:

$$R_S(f) := \frac{1}{n} \sum_{i=1}^n [f(x_i) - y_i]^2.$$

- By the **law of large numbers**, we have:

$$R_S(f) \longrightarrow R(f) \quad \text{as } n \rightarrow \infty.$$

# Example of Expected and Empirical Risk using GMM

Suppose  $(x, y)$  follows GMM, and the function  $f(x) = \theta x$ .

- The expected risk  $R(f)$  is given by

$$\begin{aligned} R(f) &= \mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(f(x), y) \\ &= \int [f(x) - y]^2 p(x, y) dx dy = \int [f(x) - y]^2 p(x|y) p(y) dx dy \\ &= \frac{1}{2} \int [f(x)]^2 p(x|y=0) dx + \frac{1}{2} \int [f(x) - 1]^2 p(x|y=1) dx \\ &= \frac{1}{2} \int [\theta x]^2 \cdot \mathcal{N}(x; \mu_1, \sigma_1^2) dx + \frac{1}{2} \int [\theta x - 1]^2 \cdot \mathcal{N}(x; \mu_2, \sigma_1^2) dx \\ &\triangleq R(\theta), \end{aligned}$$

where  $p(x, y)$  is the joint density, and  $\mathcal{N}(x; \mu, \sigma^2)$  is the Gaussian density defined by

$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}.$$

- The empirical risk  $R_S(f)$  over a training sample is given by

$$R_S(f) = \frac{1}{n} \sum_{i=1}^n [\theta x_i - y_i]^2 \triangleq R_S(\theta).$$



# Hypothesis Class

In practice, we cannot evaluate all possible functions  $f$ . Instead, we restrict our search to a family of functions called a **hypothesis class**  $\mathcal{H}$ . Each function  $h \in \mathcal{H}$  is called a **hypothesis**.

- The collection of all linear models or the collection of all two-layer neural networks:

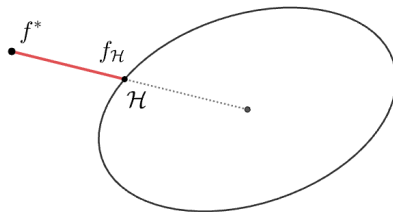
$$\mathcal{H}_1 = \{h : h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}\}$$

$$\mathcal{H}_2 = \{h : h(\mathbf{x}) = \mathbf{v}^\top \phi(\mathbf{W}\mathbf{x})\}.$$

- A learning algorithm aims to find the best hypothesis  $h \in \mathcal{H}$  that minimizes the expected risk:

$$f_{\mathcal{H}} := \operatorname{argmin}_{f \in \mathcal{H}} R(f).$$

- The difference  $\|f^* - f_{\mathcal{H}}\|$  is called the **approximation error**, where  $f^*$  is the ground true function.
- The **Universal Approximation Theorem (UAT)** implies  $\|f^* - f_{\mathcal{H}}\| \approx 0$  if  $f_{\mathcal{H}} \in \mathcal{H}_2$ .

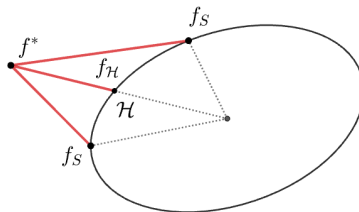


# Decomposition of Expected Risk

- Given a learned hypothesis  $f_S$  from a sample  $S$ , the expected risk of  $f_S$  can be decomposed as:

$$R(f_S) = \underbrace{R_S(f_S)}_{\text{Training Error}} + \underbrace{[R(f_S) - R_S(f_S)]}_{\text{Generalization Error}}.$$

- The **generalization error** is the difference between the expected risk and the empirical risk.



- In practice, the generalization error is estimated using the **test error** on an independent test set.

# Bounding the Generalization Error

- The generalization error can be upper bounded by the **complexity** of the hypothesis class:

$$\sup_{h \in \mathcal{H}} |R(h) - R_S(h)| \leq \text{Complexity Term},$$

where the “Complexity Term” quantifies how flexible or complex the hypothesis class  $\mathcal{H}$  is.

- One commonly used complexity measure is the **(empirical) Rademacher complexity**:

$$\mathfrak{R}_S(\mathcal{H}) := \mathbb{E}_\sigma \left[ \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), \sigma_i) \right],$$

where  $\ell(h(x), \sigma) = \sigma h(x)$  and  $\sigma_i \in \{-1, 1\}$  are i.i.d. Rademacher random variables (uniformly distributed), i.e.,  $\sigma \sim \mathcal{U}\{-1, 1\}$ , and the expectation is taken over these random labels.

- Rademacher complexity measures the ability of the hypothesis class to fit **random noise** (i.e., how well the hypothesis class can fit random labels).
- Using model complexity, we can derive the following generalization bound:

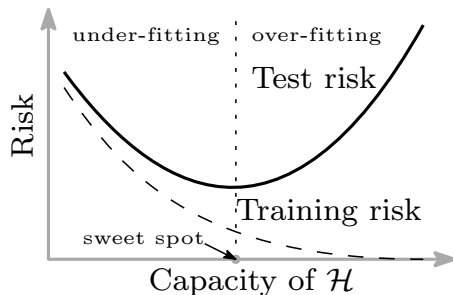
$$R(f_S) \leq R_S(f_S) + \mathfrak{R}_S(\mathcal{H}) + \tilde{O}(n^{-1}),$$

where the expected risk is upper bounded by the training error and the complexity of the model.

# Model Complexity Trade-Off

The expected risk  $R(f_S)$  is upper bounded by the training error and the model complexity:

$$R(f_S) \leq R_S(f_S) + \mathfrak{R}_S(\mathcal{H}) + \tilde{\mathcal{O}}(n^{-1}).$$



## Key Insights on Generalization Bound

- If the model is too simple, it may fail to fit the training data well. This is known as **underfitting**.
- Conversely, if the model is highly flexible, it may achieve low training error, but perform poorly on unseen data. This is known as **overfitting**.
- The goal is to find a “**sweet spot**” balancing underfitting and overfitting to minimize the overall expected risk.

Optimal Hypothesis  $f^*$ 

**Claim:**  $f^*(x) = \mathbb{E}[y|x]$  is the **optimal hypothesis** that minimizes the expected risk.

Proof.

For any function  $f$ , we can decompose the expected risk as follows:

$$\begin{aligned} R(f) &= \mathbb{E}(f - y)^2 = \mathbb{E}(f - f^* + f^* - y)^2 \\ &= \mathbb{E}(f - f^*)^2 + 2\mathbb{E}(f - f^*)(f^* - y) + \mathbb{E}(f^* - y)^2 \\ &= \mathbb{E}(f - f^*)^2 + \mathbb{E}(f^* - y)^2 \\ &\geq \mathbb{E}(f^* - y)^2 \\ &= R(f^*) \end{aligned}$$

where the cross term  $\mathbb{E}(f - f^*)(f^* - y) = 0$ , because  $f^*(x) = \mathbb{E}[y|x]$ . □

- This is another **existence** result.
- The optimal hypothesis  $f^*$  is not directly accessible unless we know the joint distribution  $\mathcal{D}$ .
- Generally, we may have  $R(f^*) \neq 0$ . For example, consider  $y = \theta x + \varepsilon$ , where  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

$$\begin{aligned} f^*(x) &= \mathbb{E}[y|x] = \mathbb{E}[\theta x + \varepsilon|x] = \theta x \\ R(f^*) &= \mathbb{E}_x[f^*(x) - y]^2 = \mathbb{E}_x[\theta x - (\theta x + \varepsilon)]^2 = \sigma^2 \implies \text{irreducible error.} \end{aligned}$$

# Bias-Variance Decomposition of Expected Risk

- The learned function  $f_S$  depends on the random sample  $S$ , making  $f_S$  a **random variable**.
- Hence, the expected risk  $R(f_S)$  is also **random**, and it varies across different random samples  $S$ .
- To capture this variability, we consider the expectation of the  $R(f_S)$  over all possible samples  $S$ , i.e.,  $\mathbb{E}_S[R(f_S)]$ .
- Let  $\bar{f} := \mathbb{E}_S[f_S]$ , the average hypothesis over all random samples  $S$ .
- Using  $\bar{f}$ , we can decompose  $\mathbb{E}_S[R(f_S)]$  as follows:

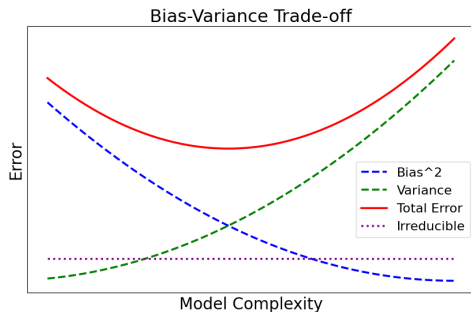
$$\begin{aligned}\mathbb{E}_S[R(f_S)] &= \mathbb{E}_S \mathbb{E}_{(x,y) \sim \mathcal{D}} [f_S(x) - y]^2 \\ &= \mathbb{E}_S \mathbb{E}_{\mathcal{D}} [f_S - f^*]^2 + R(f^*) \\ &= \mathbb{E}_S \mathbb{E}_{\mathcal{D}} [f_S - \bar{f} + \bar{f} - f^*]^2 + R(f^*) \\ &= \mathbb{E}_S \mathbb{E}_{\mathcal{D}} [(f_S - \bar{f})^2 + (\bar{f} - f^*)^2] + R(f^*) \\ &= \underbrace{\mathbb{E}_S (f_S - \bar{f})^2}_{\text{Variance term}} + \underbrace{\mathbb{E}_{\mathcal{D}} (\bar{f} - f^*)^2}_{\text{Bias term}} + \underbrace{R(f^*)}_{\text{irreducible}}\end{aligned}$$

where the cross term  $\mathbb{E}_{S,\mathcal{D}}(f_S - \mathbb{E}_S[f_S])(\mathbb{E}_S[f_S] - f^*) = 0$  cancels out.

# Bias-Variance Trade-Off

The expected risk  $\mathbb{E}_S[R(f_S)]$  can be broken down into three parts:

- **Squared Bias:**  $\mathbb{E}_{\mathcal{D}}[(f^* - \mathbb{E}_S(f_S))^2]$  measures the error from approximating the optimal function  $f^*$  with the learned model  $f_S$ . It reflects the error caused by using a simple model that cannot capture all the data patterns.
- **Variance:**  $\text{Var}(f_S) = \mathbb{E}_S[(f_S - \mathbb{E}_S(f_S))^2]$  measures how much the learned function  $f_S$  varies with different training samples. It represents the error due to the model's sensitivity to fluctuations in the random training sample  $S$ .
- **Irreducible Error:**  $R(f^*)$  represents the inherent noise in the data, which no model can eliminate. It is the error we cannot reduce.



- **High bias, low variance:** Simple models (e.g., linear models) have *low variance* since they are less sensitive to training data, but have *high bias* because they are too simple to capture all patterns in the data.
- **Low bias, high variance:** Complex models (e.g., polynomial model) have *low bias* as they can model complex relations, but *high variance* due to overfitting to the training data.

# Summary of Statistical Learning Theory

- The goal is to find a hypothesis  $f$  within a hypothesis class  $\mathcal{H}$  that minimizes the **expected risk**:

$$R(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [(f(x) - y)^2].$$

- Since the underlying distribution  $\mathcal{D}$  is **unknown**, we approximate  $f$  by minimizing the **empirical risk** based on a **random** training sample  $S$ :

$$R_S(f) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2.$$

- Using model complexity  $\mathfrak{R}_S(\mathcal{H})$ , the expected risk is upper bound as:

$$R(f_S) \leq R_S(f_S) + \mathfrak{R}_S(\mathcal{H}) + \tilde{O}(n^{-1}),$$

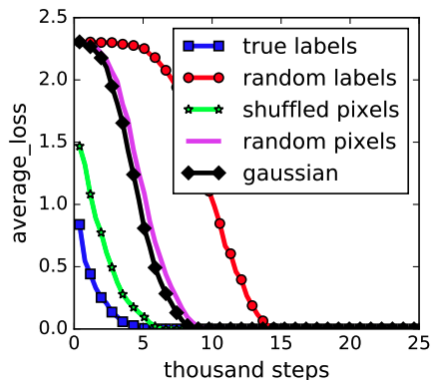
- By considering variations across different random training samples  $S$ , the expected risk  $\mathbb{E}_S[R(f_S)]$  can be decomposed into three components: **bias**, **variance**, and **irreducible error**:
  - **High bias, low variance**: Simple models **underfit** and miss important patterns in the data.
  - **Low bias, high variance**: Complex models **overfit** and perform poorly on unseen data.
- Find the “sweet spot” between underfitting and overfitting to minimize the overall expected risk.



# Outline

- 1 Statistical Learning Theory
- 2 Regularization
- 3 Hyperparameter Tune
- 4 Overparameterization

# DNNs Can Fit Random Labels and Random Data



- **Label corruption:** Replace true label with *random label*
- **Shuffled pixels:** The pixels of each image are rearranged using a *fixed* random permutation
- **Random pixels:** Each image has a *unique* random arrangement of pixels).
- **Gaussian:** The pixels in images are replaced with random Gaussian *noise*.
- **Average loss:** *Training error* using the cross-entropy loss

## Key Observation

DNNs can perfectly fit random labels or data, achieving **zero training error** even on completely unstructured inputs.

# Weight Decay

- Regularization typically involves adding an extra term, called the **regularizer**, to the training loss:

$$\mathcal{L}_\lambda(\boldsymbol{\theta}) := \mathcal{L}(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2,$$

where  $\lambda > 0$  is the **regularization hyperparameter**, and  $\|\cdot\|$  is the Euclidean norm.

- In deep learning, this regularization is known as **weight decay** because gradient descent on the regularized loss automatically shrinks (or decays) parameter  $\boldsymbol{\theta}$  by the factor  $(1 - \eta\lambda)$ :

$$\begin{aligned} \boldsymbol{\theta}^+ &= \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}_\lambda(\boldsymbol{\theta}) = \boldsymbol{\theta} - \eta [\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}] \\ &= \underbrace{(1 - \eta\lambda)}_{\text{decaying weights}} \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}). \end{aligned}$$

- However,  $\boldsymbol{\theta}$  does **not** shrink to zero, as it must maintain a certain value to minimize the cost  $\mathcal{L}(\boldsymbol{\theta})$ .

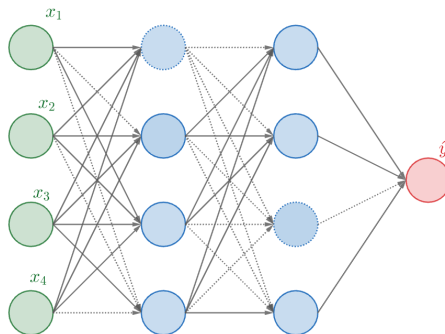
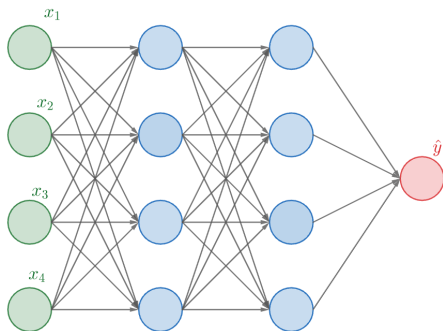
# Interpretation: Sparsity

- The regularized optimization can be reformulated as:

$$\min_{\theta} \mathcal{L}(\theta), \quad \text{s.t.} \quad \|\theta\| \leq C_{\lambda},$$

where  $C_{\lambda} > 0$  is a constant that depends on  $\lambda$ .

- In deep learning,  $\theta$  is called **sparse** if most parameters are zero or close to zero (i.e.,  $\theta_i \approx 0$ ).
- Sparse  $\theta$  reduces the flexibility and complexity of the DNN, leading to a **simpler** model.

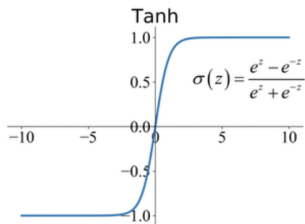


# Interpretation: Linearity

Consider a simple two-layer neural network:

$$f_{\theta}(x) = \sum_{i=1}^n v_i \phi(w_i x),$$

where  $x \in \mathbb{R}$  is a scalar and  $\phi(\cdot)$  is tanh.



- When  $w_i \approx 0$ , then  $w_i x \approx 0$ , and the network operates near the **linear** region of tanh:

$$v_i \phi(w_i x) \approx v_i (w_i x) \approx (v_i w_i) x = u_i x \implies \text{a linear model,}$$

where  $u_i := v_i w_i$ .

- If  $v_i \approx 0$ , then

$$v_i \phi(w_i x) \approx 0,$$

indicating **fewer** neurons are used.

# Interpretation: Stability

A learning algorithm is **stable** if small changes to its input do not result in large changes to its output.

- Consider the same two-layer neural network:

$$f_{\theta}(x) = \sum_{i=1}^n v_i \phi(w_i x).$$

- The derivative of  $f_{\theta}$  with respect to the **input**  $x$  is:

$$\nabla_x f_{\theta}(x) = \sum_{i=1}^n v_i \phi'(w_i x) w_i.$$

- If either  $v_i$  or  $w_i$  is small, then  $\nabla_x f_{\theta}(x)$  is small.
- Hence, DNNs with *sparse* parameters are generally more **stable** than those with dense parameters.

# Dropout Regularization

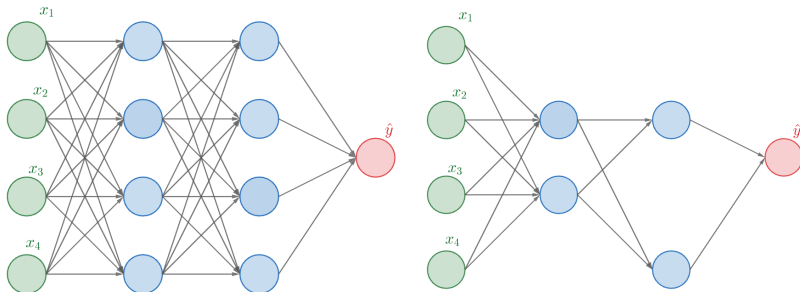
Recall the forward propagation:

$$\mathbf{z}^\ell = \mathbf{W}^\ell \mathbf{x}^{\ell-1}, \quad \mathbf{x}^\ell = \phi(\mathbf{z}^\ell).$$

- During **training**, each neuron is randomly **dropped** with probability  $p$  (a **hyperparameter**):

$$\hat{\mathbf{x}}^{\ell-1} = \mathbf{r}^\ell \odot \mathbf{x}^{\ell-1}, \quad \mathbf{z}^\ell = \mathbf{W}^\ell \hat{\mathbf{x}}^{\ell-1}$$

where  $\mathbf{r}_i^\ell \sim \text{Bernoulli}(p)$ .

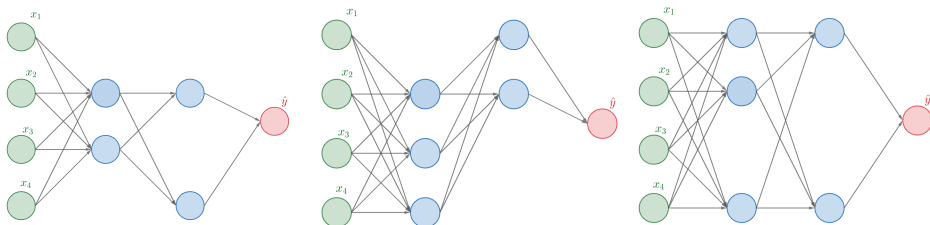


- The gradient update applies only to a **thinned subnetwork** of the network.
- At **test time**, dropout is turned off, and weights are scaled by  $p$  to respect the dropout probability:

$$\mathbf{z}^\ell = p \mathbf{W}^\ell \mathbf{x}^{\ell-1}.$$

# Interpretation: Implicit Ensemble Learning

- By randomly dropping units, a different **thinned subnet** is trained at each gradient descent step.
- With  $n$  neurons in the full network, we are effectively training  $2^n$  **different subnets** simultaneously that all **share** the same weights.
- At test time, the output is an **ensemble** prediction, aggregating the contributions of all subnets.



## Key Insight

Dropout ensures that **no** single neuron or small group of neurons can dominate the prediction. By **spreading** the responsibility across all units, it improves model robustness to the input change and prevents overfitting.



# Summary

- DNNs can fit random labels and data, achieving **zero training error**.
- Weight decay controls large weights, promoting **sparsity**, **linearity**, and **stability**.
- During training, dropout randomly drops units, effectively training an **exponential number of thinned subnets** simultaneously.
- At test time, the output is an **ensemble** prediction, aggregating contributions from all subnets.

# Outline

- 1 Statistical Learning Theory
- 2 Regularization
- 3 **Hyperparameter Tune**
- 4 Overparameterization

# Validation Set

- Split the dataset into three parts: **training set**, **validation set**, and **test set**.
- Build the model using the *training set*.
- Optimize or tune hyperparameters on the *validation set*.
- After tuning, **evaluate** the final model on the test set.
- Suggested split ratios:
  - For datasets between 100 and 1,000,000 samples: 60/20/20.
  - For datasets larger than 1,000,000 samples: 98/1/1.
- Ensure the validation and test sets come from the **same** distribution.
  - Example: Training and validation images from the web, but test images from user cell phones can cause a mismatch.

# Tuning Process

- **Random Search**: Randomly sample hyperparameters; more efficient than grid search when some hyperparameters are less important.
- **Hyperband/Successive Halving**: Dynamically allocate resources and discard poor configurations early, ideal for deep networks with long training times.
- Start with **coarse** tuning, then **refine** gradually.
- Use **log scale** for hyperparameter search when appropriate.
- Leverage **parallelization** to run multiple experiments simultaneously to accelerate the search.

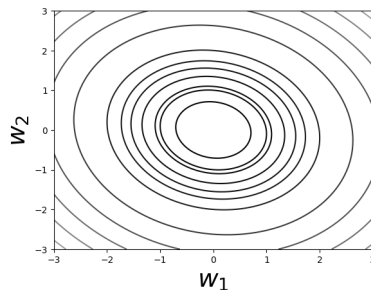
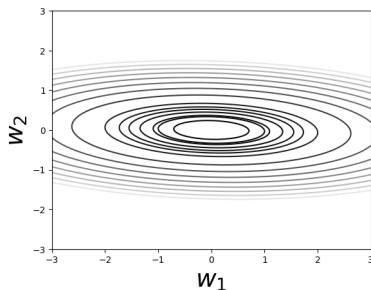
# Input Normalization

- Normalize the inputs using **training set**:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{x}_i = x_i - \mu, \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n \bar{x}_i^2, \quad \hat{x}_i = \bar{x}_i / \sigma,$$

where all operations are taken element-wise.

- Consider a binary classification problem using linear model:  $f_{\theta}(x) = w_1 x_1 + w_2 x_2$ 
  - if  $x_1 = \mathcal{O}(100)$  and  $x_2 = \mathcal{O}(1)$ , to have output  $f_{\theta} = \mathcal{O}(1)$ , we must have  $w_1 = \mathcal{O}(\frac{1}{100})$  and  $w_2 = \mathcal{O}(1)$ .
  - After normalization,  $\bar{x}_1 = \mathcal{O}(1)$  and  $\bar{x}_2 = \mathcal{O}(1)$ , so we have  $w_1 = \mathcal{O}(1)$  and  $w_2 = \mathcal{O}(1)$ .



- At test time, apply  $\mu$  and  $\sigma$  from **training** to test set.

# Learning Rate Decay

- Recall that an **epoch**  $k$  is one pass through all **mini-batches** in SGD
- Instead of using a fixed learning rate, one can consider using **learning rate decay**

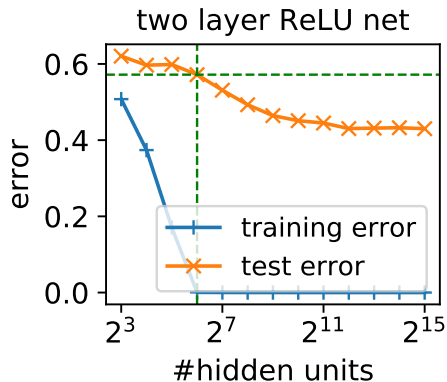
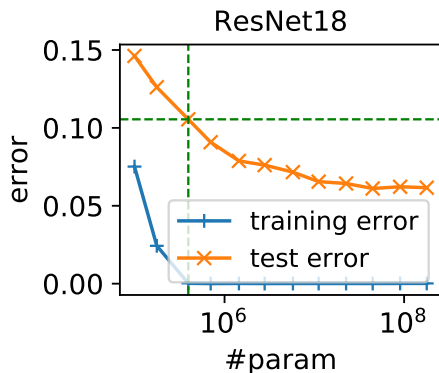
$$\eta_k = \frac{\eta}{k}, \quad \eta_k = \frac{\eta}{\sqrt{k}}, \quad \eta_k = (0.95)^k \eta$$

# Outline

- 1 Statistical Learning Theory
- 2 Regularization
- 3 Hyperparameter Tune
- 4 Overparameterization

# Overparameterization

- A deep neural network (DNN) is said to be **overparameterized** when the number of neurons or parameters is much larger than the number of training samples.
- This might seem counterintuitive, but it has been found to be surprisingly beneficial in practice.





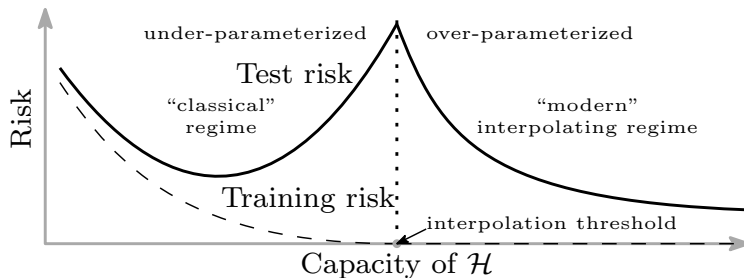
# Double Descent

Overparameterized neural networks can **perfectly fit** or **interpolate** the training data.

- Mathematically, there exists a set of parameters  $\theta$  such that

$$f_{\theta}(x_i) = y_i, \quad \forall i \in [n]. \quad (1)$$

- Overparameterization implies there are **infinitely** many interpolation solutions.
- Some interpolation solutions generalize much better than those in the *underparameterized* regime. This phenomenon is called **double descent**.



# Implicit Regularization

- It is important to understand that different global minima lead to varying **test** performances.
- A **flat** minimum typically results in better generalization than a **sharp** minimum.
- Different optimizers may converge to different minima, each with different generalization outcomes. This is known as **implicit regularization**.
- Thus, even if your current optimizer achieves low training error, tuning or adjusting it may still be necessary to achieve better test performance.

