

Learning Theory & Regularization

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

Machine Learning

Outline

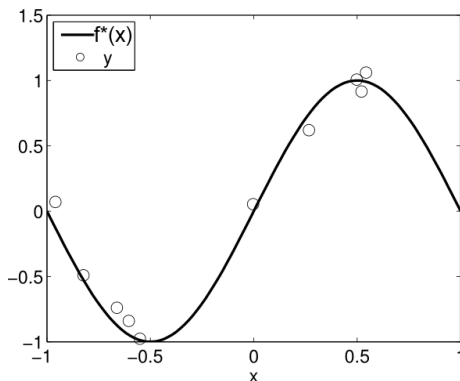
- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization
 - Weight Decay
 - Validation

Outline

- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization
 - Weight Decay
 - Validation

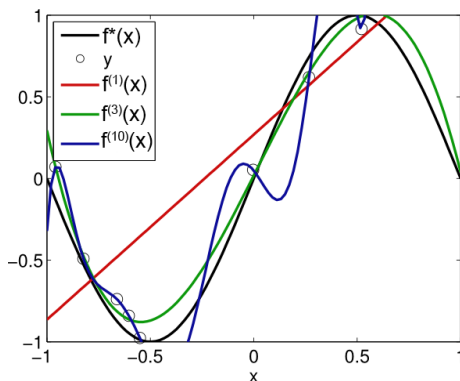
Which Polynomial Degree Is Better? I

- Given a training set $\mathbb{X} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$ i.i.d. sampled from $P(\mathbf{x}, \mathbf{y})$
- Assume $P(\mathbf{x}, \mathbf{y}) = P(\mathbf{y} | \mathbf{x})P(\mathbf{x})$, where
 - $P(\mathbf{x}) \sim \text{Uniform}(-1, 1)$
 - $\mathbf{y} = \sin(\pi \mathbf{x}) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$



Which Polynomial Degree Is Better? II

- Consider 3 unregularized polynomial regressors of degrees $P = 1, 3,$ and 10



- Which one would you pick? Probably not $P = 1$ nor $P = 10$
- Note that $P = 10$ has **zero** training error
 - Any N points can be perfectly fitted by a polynomial of degree $N - 1$

Empirical Error vs. Generalization Error

- In ML, we usually “learn” a function by minimizing the **empirical error/risk** defined over a training set of size N :

$$C_N(\mathbf{w}) \text{ or } C_N[f] = \frac{1}{N} \sum_{i=1}^N \text{loss} \left(f(\mathbf{x}^{(i)}; \mathbf{w}), \mathbf{y}^{(i)} \right)$$

- E.g., $C_N(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \left(y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \right)^2$ in linear regression
- But our goal is to have a low **generalization error/risk** defined over the underlying data distribution:

$$C(\mathbf{w}) \text{ or } C[f] = \int \text{loss} (f(\mathbf{x}; \mathbf{w}), y) dP(\mathbf{x}, y)$$

- Can be estimated by the **testing error**
 $C_{N'}(\mathbf{w}) = \frac{1}{N'} \sum_{i=1}^{N'} \text{loss} \left(f(\mathbf{x}'^{(i)}; \mathbf{w}), \mathbf{y}'^{(i)} \right)$ defined over the testing set
 $\mathbb{X}' = \{(\mathbf{x}'^{(i)}, \mathbf{y}'^{(i)})\}_{i=1}^{N'}$

- Does a low $C_N[f]$ implies low $C[f]$? No, as $P = 10$ indicates

No-Free-Lunch Theorem

- Why $C[f]$ is defined over a *particular* data generating distribution P ?

Theorem (No-Free-Lunch Theorem [4])

Averaged over all possible data generating distributions, every classification algorithm has the same error rate when classifying unseen points.

- No machine learning algorithm is better than any other universally
- The goal of ML is *not* to seek a universally good learning algorithm
- Instead, a good algorithm that performs well on data drawn from a *particular P* we care about

Learning Theory

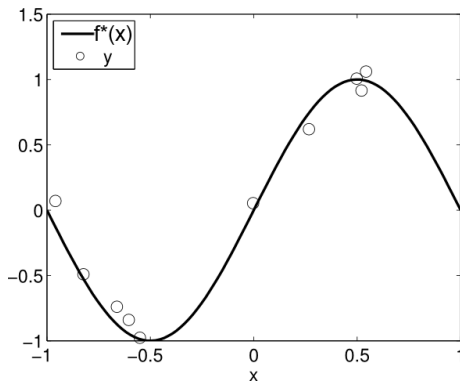
- Let $f^* = \arg \min_f C[f]$ be the best possible function we can get
- Since we are seeking a prediction function in a model (hypothesis space) \mathbb{F} , this is what can have at best: $f_{\mathbb{F}}^* = \arg \min_{f \in \mathbb{F}} C[f]$
- But we only minimizes empirical errors on limited examples of size N , this is what we actually have $f_N = \arg \min_{f \in \mathbb{F}} C_N[f]$
 - Ignoring numerical errors (due to, e.g., numerical optimization)
- **Learning theory**: how to characterize

$$C[f_N] = \int \text{loss}(f_N(\mathbf{x}; \mathbf{w}), y) d\mathbf{P}(\mathbf{x}, y)?$$

- Not to confuse $C[f_N]$ with $C_N[f]$
- Bounding methods
- Decomposition methods

Bounding Methods I

- $\min_f C[f] = C[f^*]$ is called the *Bayes error*
 - Larger than 0 when there is randomness in $P(y|x)$
 - E.g., in our regression problem: $y = f^*(x; \mathbf{w}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
- *Cannot be avoided* even we know $P(x, y)$ in the ground truth



- So, our target is to make $C[f_N]$ as close to $C[f^*]$ as possible

Bounding Methods II

- Let $\mathcal{E} = C[f_N] - C[f^*]$ be the *excess error*
- We have

$$\mathcal{E} = \underbrace{C[f_{\mathbb{F}}^*] - C[f^*]}_{\mathcal{E}_{\text{app}}} + \underbrace{C[f_N] - C[f_{\mathbb{F}}^*]}_{\mathcal{E}_{\text{est}}}$$

- \mathcal{E}_{app} is called the *approximation error*
- \mathcal{E}_{est} is called the *estimation error*
- How to reduce these errors?
- We can reduce \mathcal{E}_{app} by choosing a *more complex* \mathbb{F}
 - A complex \mathbb{F} has a larger capacity
 - E.g., larger polynomial degree P in polynomial regression
- How to reduce \mathcal{E}_{est} ?

Bounding Methods III

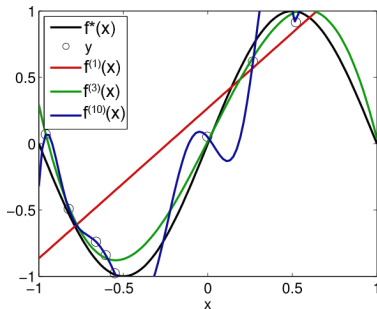
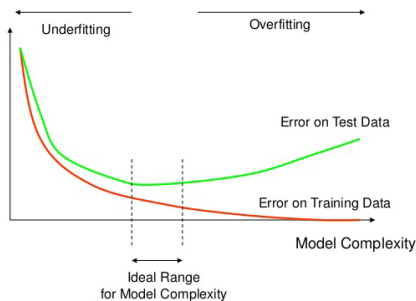
- Bounds of \mathcal{E}_{est} for, e.g., binary classifiers [1, 2, 3]:

$$\mathcal{E}_{\text{est}} = O \left[\left(\frac{\text{Complexity}(\mathbb{F}) \log N}{N} \right)^\alpha \right], \alpha \in \left[\frac{1}{2}, 1 \right], \text{ with high probability}$$

- So, to reduce \mathcal{E}_{est} , we should either have
 - **Simpler model** (e.g., smaller polynomial degree P), or
 - Larger training set

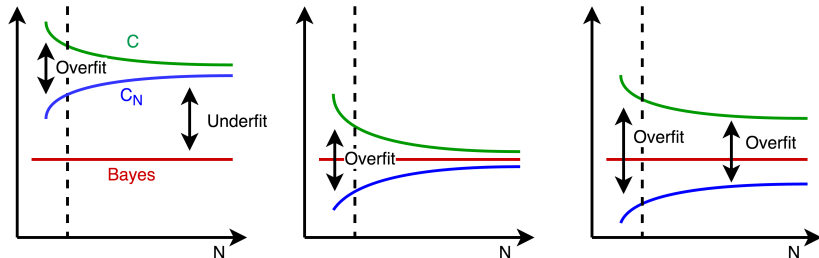
Model Complexity, Overfit, and Underfit

- Too simple a model leads to high \mathcal{E}_{app} due to **underfitting**
 - f_N fails to capture the shape of f^*
 - **High** training error; **high** testing error (given a sufficiently large N)
- Too complex a model leads to high \mathcal{E}_{est} due to **overfitting**
 - f_N captures not only the shape of f^* but also some spurious patterns (e.g., noise) local to a particular training set
 - **Low** training error; **high** testing error



Sample Complexity and Learning Curves

- How many training examples (N) are sufficient?
- Different models/algorithms may have different *sample complexity*
 - I.e., the N required to learn a target function with specified generalizability
- Can be visualized using the *learning curves*
- Too small N results in overfit regardless of model complexity



Decomposition Methods

- Bounding methods analyze $C[f_N]$ qualitatively
 - General, as no (or weak) assumption on data distribution is made
- However, in practice, these bounds are too loose to quantify $C[f_N]$
- In some particular situations, we can *decompose* $C[f_N]$ into multiple meaningful terms
- Assume particular
 - Loss function $\text{loss}(\cdot)$, and
 - Data generating distribution $P(x, y)$
- Require knowledge about the *point estimation*

Outline

- ① Learning Theory
- ② **Point Estimation: Bias and Variance**
 - Consistency*
- ③ Decomposing Generalization Error
- ④ Regularization
 - Weight Decay
 - Validation

Sample Mean and Variance

- **Point estimation** is the attempt to estimate some fixed but unknown quantity θ of a random variable by using sample data
- Let $\mathbb{X} = \{x^{(1)}, \dots, x^{(n)}\}$ be a set of n i.i.d. samples of a random variable x , a **point estimator** or **statistic** is a function of the data:

$$\hat{\theta}_n = g(x^{(1)}, \dots, x^{(n)})$$

- The value $\hat{\theta}_n$ is called the **estimate** of θ
- Sample mean: $\hat{\mu}_x = \frac{1}{n} \sum_i x^{(i)}$
- Sample variance: $\hat{\sigma}_x = \frac{1}{n} \sum_i (x^{(i)} - \hat{\mu}_x)^2$
- How good are these estimators?

Bias & Variance

- **Bias** of an estimator:

$$\text{bias}(\hat{\theta}_n) = E_{\mathbb{X}}(\hat{\theta}_n) - \theta$$

- Here, the expectation is defined over **all possible \mathbb{X} 's of size n** , i.e.,
 $E_{\mathbb{X}}(\hat{\theta}_n) = \int \hat{\theta}_n dP(\mathbb{X})$
- We call a statistic **unbiased estimator** iff it has zero bias
- **Variance** of an estimator:

$$\text{Var}_{\mathbb{X}}(\hat{\theta}_n) = E_{\mathbb{X}} \left[(\hat{\theta}_n - E_{\mathbb{X}}[\hat{\theta}_n])^2 \right]$$

- Is $\hat{\mu}_x = \frac{1}{n} \sum_i x^{(i)}$ an unbiased estimator of μ_x ? Yes [Homework]
- What much is $\text{Var}_{\mathbb{X}}(\hat{\mu}_x)$?

Variance of $\hat{\mu}_x$

$$\begin{aligned}\text{Var}_{\mathbb{X}}(\hat{\mu}) &= \mathbb{E}_{\mathbb{X}}[(\hat{\mu} - \mathbb{E}_{\mathbb{X}}[\hat{\mu}])^2] = \mathbb{E}[\hat{\mu}^2 - 2\hat{\mu}\mu + \mu^2] = \mathbb{E}[\hat{\mu}^2] - \mu^2 \\&= \mathbb{E}\left[\frac{1}{n^2} \sum_{i,j} x^{(i)} x^{(j)}\right] - \mu^2 = \frac{1}{n^2} \sum_{i,j} \mathbb{E}[x^{(i)} x^{(j)}] - \mu^2 \\&= \frac{1}{n^2} \left(\sum_{i=j} \mathbb{E}[x^{(i)} x^{(j)}] + \sum_{i \neq j} \mathbb{E}[x^{(i)} x^{(j)}] \right) - \mu^2 \\&= \frac{1}{n^2} \left(\sum_i \mathbb{E}[x^{(i)2}] + n(n-1) \mathbb{E}[x^{(i)}] \mathbb{E}[x^{(j)}] \right) - \mu^2 \\&= \frac{1}{n} \mathbb{E}[x^2] + \frac{(n-1)}{n} \mu^2 - \mu^2 = \frac{1}{n} (\mathbb{E}[x^2] - \mu^2) = \frac{1}{n} \sigma_x^2\end{aligned}$$

- The variance of $\hat{\mu}_x$ diminishes as $n \rightarrow \infty$

Unbiased Estimator of σ_x

- Is $\hat{\sigma}_x = \frac{1}{n} \sum_i (x^{(i)} - \hat{\mu}_x)^2$ and an unbiased estimator of σ_x ? **No**

$$\begin{aligned} E_{\mathbb{X}}[\hat{\sigma}] &= E\left[\frac{1}{n} \sum_i (x^{(i)} - \hat{\mu})^2\right] = E\left[\frac{1}{n} (\sum_i x^{(i)2} - 2 \sum_i x^{(i)} \hat{\mu} + \sum_i \hat{\mu}^2)\right] \\ &= E\left[\frac{1}{n} (\sum_i x^{(i)2} - n \hat{\mu}^2)\right] = \frac{1}{n} (\sum_i E[x^{(i)2}] - n E[\hat{\mu}^2]) \\ &= E[x^2] - E[\hat{\mu}^2] = E[(x - \mu)^2 + 2x\mu - \mu^2] - E[\hat{\mu}^2] \\ &= (\sigma^2 + \mu^2) - (E[\hat{\mu}^2]) \\ &= \sigma^2 + \mu^2 - \frac{1}{n} \sigma^2 - \mu^2 = \frac{n-1}{n} \sigma^2 \neq \sigma^2 \end{aligned}$$

- What's the unbiased estimator of σ_x ?

$$\hat{\sigma}_x = \frac{n}{n-1} \left(\frac{1}{n} \sum_i (x^{(i)} - \hat{\mu}_x)^2 \right) = \frac{1}{n-1} \sum_i (x^{(i)} - \hat{\mu}_x)^2$$

Mean Square Error

- *Mean square error* of an estimator:

$$\text{MSE}(\hat{\theta}_n) = \mathbb{E}_{\mathbb{X}} [(\hat{\theta}_n - \theta)^2]$$

- Can be decomposed into the bias and variance:

$$\begin{aligned}\mathbb{E}_{\mathbb{X}} [(\hat{\theta}_n - \theta)^2] &= \mathbb{E} [(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n] - \mathbb{E}[\hat{\theta}_n] + \theta)^2] \\&= \mathbb{E} [(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2 + (\mathbb{E}[\hat{\theta}_n] - \theta)^2 + 2(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])(\mathbb{E}[\hat{\theta}_n] - \theta)] \\&= \mathbb{E} [(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2] + \mathbb{E} [(\mathbb{E}[\hat{\theta}_n] - \theta)^2] + 2\mathbb{E}(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])(\mathbb{E}[\hat{\theta}_n] - \theta) \\&= \mathbb{E} [(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2] + (\mathbb{E}[\hat{\theta}_n] - \theta)^2 + 2 \cdot 0 \cdot (\mathbb{E}[\hat{\theta}_n] - \theta) \\&= \text{Var}_{\mathbb{X}}(\hat{\theta}_n) + \text{bias}(\hat{\theta}_n)^2\end{aligned}$$

- MSE of an unbiased estimator is its variance

Outline

- 1 Learning Theory
- 2 Point Estimation: Bias and Variance**
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization
 - Weight Decay
 - Validation

Consistency

- So far, we discussed the “goodness” of an estimator based on samples of fixed size
- If we have more samples, will the estimate become more accurate?
- An estimator is (weak) *consistent* iff:

$$\lim_{n \rightarrow \infty} \hat{\theta}_n \xrightarrow{\text{Pr}} \theta,$$

where $\xrightarrow{\text{Pr}}$ means “converge in probability”

- Strong consistent iff “converge almost surely”

Law of Large Numbers

Theorem (Weak Law of Large Numbers)

The sample mean $\hat{\mu}_x = \frac{1}{n} \sum_i x^{(i)}$ is a consistent estimator of μ_x , i.e., $\lim_{n \rightarrow \infty} \Pr(|\hat{\mu}_{x,n} - \mu_x| < \varepsilon) = 1$ for any $\varepsilon > 0$.

Theorem (Strong Law of Large Numbers)

In addition, $\hat{\mu}_x$ is a strong consistent estimator: $\Pr(\lim_{n \rightarrow \infty} \hat{\mu}_{x,n} = \mu_x) = 1$.

Outline

- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error**
- 4 Regularization
 - Weight Decay
 - Validation

Expected Generalization Error

- In ML, we get $f_N = \arg \min_{f \in \mathbb{F}} C_N[f]$ by minimizing the empirical error over a training set of size N
- How to decompose the generalization error $C[f_N]$?
- Regard $f_N(\mathbf{x})$ as an estimate of true label y given \mathbf{x}
 - f_N an estimator mapped from i.i.d. samples in the training set \mathbb{X}
- To evaluate the estimator f_N , we consider the expected generalization error:

$$\begin{aligned} E_{\mathbb{X}}(C[f_N]) &= E_{\mathbb{X}}[\int \text{loss}(f_N(\mathbf{x}) - y) d\mathbf{P}(\mathbf{x}, y)] \\ &= E_{\mathbb{X}, \mathbf{x}, y}[\text{loss}(f_N(\mathbf{x}) - y)] \\ &= E_{\mathbf{x}}(E_{\mathbb{X}, y}[\text{loss}(f_N(\mathbf{x}) - y) | \mathbf{x} = \mathbf{x}]) \end{aligned}$$

- There's a simple decomposition of $E_{\mathbb{X}, y}[\text{loss}(f_N(\mathbf{x}) - y) | \mathbf{x}]$ for linear/polynomial regression

Example: Linear/Polynomial Regression

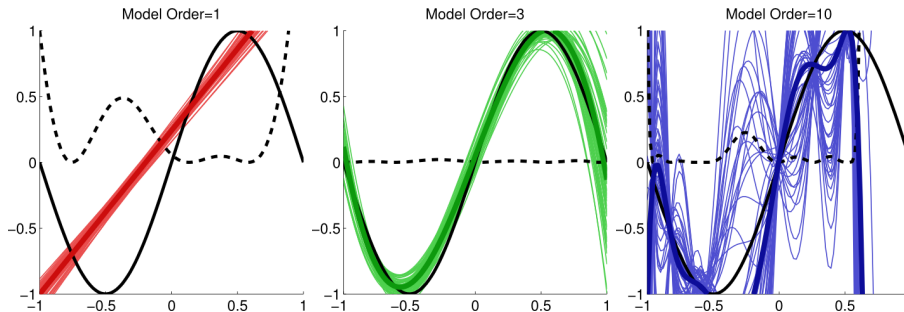
- In linear/polynomial regression, we have
 - $\text{loss}(\cdot) = (\cdot)^2$ a squared loss
 - $y = f^*(\mathbf{x}) + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, thus $\mathbb{E}_y[y|\mathbf{x}] = f^*(\mathbf{x})$ and $\text{Var}_y[y|\mathbf{x}] = \sigma^2$
- We can decompose the mean square error:

$$\begin{aligned}\mathbb{E}_{\mathbb{X},y} [\text{loss}(f_N(\mathbf{x}) - y)|\mathbf{x}] &= \mathbb{E}_{\mathbb{X},y} [(f_N(\mathbf{x}) - y)^2|\mathbf{x}] \\&= \mathbb{E}_{\mathbb{X},y} [y^2 + f_N(\mathbf{x})^2 - 2f_N(\mathbf{x})y|\mathbf{x}] \\&= \mathbb{E}_y[y^2|\mathbf{x}] + \mathbb{E}_{\mathbb{X}}[f_N(\mathbf{x})^2|\mathbf{x}] - 2\mathbb{E}_{\mathbb{X},y}[f_N(\mathbf{x})y|\mathbf{x}] \\&= (\text{Var}_y[y|\mathbf{x}] + \mathbb{E}_y[y|\mathbf{x}]^2) + (\text{Var}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] + \mathbb{E}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}]^2) \\&\quad - 2\mathbb{E}_y[y|\mathbf{x}]\mathbb{E}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] \\&= \text{Var}_y[y|\mathbf{x}] + \text{Var}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] + (\mathbb{E}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] - \mathbb{E}_y[y|\mathbf{x}])^2 \\&= \text{Var}_y[y|\mathbf{x}] + \text{Var}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] + \mathbb{E}_{\mathbb{X}}[f_N(\mathbf{x}) - f^*(\mathbf{x})|\mathbf{x}]^2 \\&= \sigma^2 + \text{Var}_{\mathbb{X}}[f_N(\mathbf{x})|\mathbf{x}] + \text{bias}[f_N(\mathbf{x})|\mathbf{x}]^2\end{aligned}$$

Bias-Variance Tradeoff I

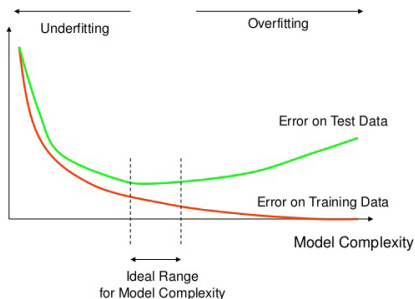
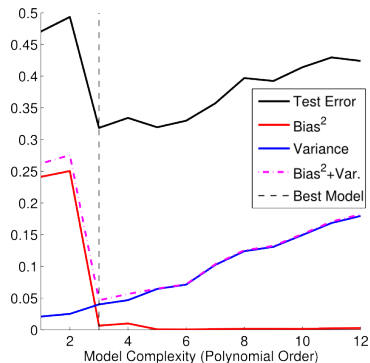
$$\begin{aligned} \mathbb{E}_{\mathbf{X}} (C[f_N]) &= \mathbb{E}_{\mathbf{x}} \left(\mathbb{E}_{\mathbf{X}, y} [\text{loss}(f_N(\mathbf{x}) - y) | \mathbf{x}] \right) \\ &= \mathbb{E}_{\mathbf{x}} \left(\sigma^2 + \text{Var}_{\mathbf{X}}[f_N(\mathbf{x}) | \mathbf{x}] + \text{bias}[f_N(\mathbf{x}) | \mathbf{x}]^2 \right) \end{aligned}$$

- The first term cannot be avoided when $P(y|\mathbf{x})$ is stochastic
- **Model complexity** controls the tradeoff between variance and bias
- E.g., polynomial regressors (dotted line = average training error):



Bias-Variance Tradeoff II

- Provides another way to understand the generalization/testing error
- Too simple a model leads to high bias or underfitting
 - **High** training error; **high** testing error (given a sufficiently large N)
- Too complex a model leads to high variance or overfitting
 - **Low** training error; **high** testing error



Outline

- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization**
 - Weight Decay
 - Validation

Regularization

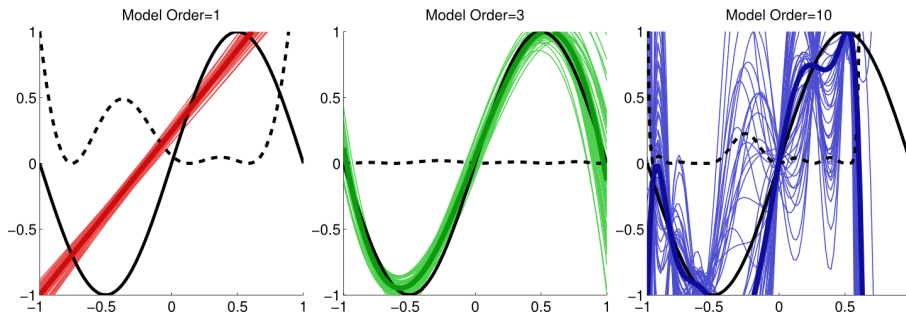
- We get $f_N = \arg \min_{f \in \mathbb{F}} C_N[f]$ by minimizing the empirical error
- But what we really care about is the generalization error $C[f_N]$
- **Regularization** refers to any technique designed to improve the generalizability of f_N
- Any idea inspired by the learning theory?
- Regularization in the cost function: **weight decay**
- Regularization during the training process: **validation**

Outline

- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization**
 - Weight Decay
 - Validation

Panelizing Complex Functions

- *Occam's razor*: among equal-performing models, the simplest one should be selected
- Idea: to add a term in the cost function that panelizes complex functions
- So, with sufficiently complex \mathbb{F} :
 - Minimizing the empirical error term reduces bias
 - Minimizing the penalty term reduces variance



What to Panelize?

- What impacts $\text{Complexity}(\mathbb{F})$ in a model?
- Some constants in the model \mathbb{F}
 - E.g., degree P in polynomial regression
- Restricts the capacity of \mathbb{F}
- However, cannot be penalized in a cost function since fixed
- Alternatively, *function parameters*
 - E.g., the parameter \mathbf{w} of a function $f(\cdot; \mathbf{w}) \in \mathbb{F}$
- Also restricts the capacity of \mathbb{F}
- Can be penalized
- But which \mathbf{w} implies a complex model?

Weight Decay

- In practice, $\mathbf{w} = \mathbf{0}$ is usually the “simplest” function
 - E.g, in binary classification for labels $\{-1, 1\}$, a perceptron with $\mathbf{w} = \mathbf{0}$ means random guessing
- **Weight decay**: to penalize the **norm** of \mathbf{w} , which is nonnegative and equals to 0 when $\mathbf{w} = \mathbf{0}$
- E.g., the **Ridge regression**:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{y} - (\mathbf{X}\mathbf{w} - b\mathbf{1})\|^2 \text{ subject to } \|\mathbf{w}\|^2 \leq T$$

for some constant $T > 0$

- In practice, we usually solve a simpler problem:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2N} \|\mathbf{y} - (\mathbf{X}\mathbf{w} - b\mathbf{1})\|^2 + \frac{\alpha}{2} \|\mathbf{w}\|^2$$

where $\alpha > 0$ is a constant representing both T and the KKT multiplier

- What does a larger α means? We prefer a more simple function

Flat Regressors

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} (\|\mathbf{y} - (X\mathbf{w} - b\mathbf{1})\|^2 + \alpha \|\mathbf{w}\|^2)$$

- The bias b is **not** regularized, why?
- We want the simplest function with $\mathbf{w} = \mathbf{0}$ means “a dummy regressor by averaging”
 - Remember R^2 (coefficient of determination)?
- However, the label y 's may not be standardized to have zero mean
- This explains why we prefer a “flat” hyperplane in the previous lecture
- We have discussed how to solve the Ridge regression problem

Sparse Weight Decay

- Alternatively we can minimize the L^1 -norm in weight decay
- E.g., **LASSO** (least absolute shrinkage and selection operator):

$$\arg \min_{\mathbf{w}, b} \frac{1}{2N} \|\mathbf{y} - (X\mathbf{w} - b\mathbf{1})\|^2 + \alpha \|\mathbf{w}\|_1$$

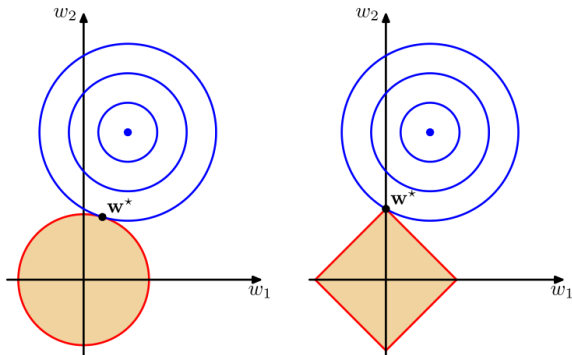
for some constant $\alpha > 0$

- Usually results in **sparse \mathbf{w}** that has many zero attributes
- Why?

Sparsity

$$\arg \min_{w,b} \frac{1}{2N} \|\mathbf{y} - (X\mathbf{w} - b\mathbf{1})\|^2 + \alpha \|\mathbf{w}\|_1$$

- The surface of the cost function is the sum of SSE (blue contours) and 1-norm (red contours)
- Optimal point locates on some axes



Elastic Net**

- LASSO can be used as a feature selection technique
 - The sparse \mathbf{w} selects explanatory variables that are most correlated to the target variable
- Limitations:
 - ① Selects at most N variables if $D > N$
 - ② No **group selection**
 - Important in some applications, e.g., gene selection problems
- **Elastic net** combines Ridge and LASSO:

$$\arg \min_{\mathbf{w}, b} \frac{1}{2N} \|\mathbf{y} - (\mathbf{X}\mathbf{w} - b\mathbf{1})\|^2 + \alpha \left(\beta \|\mathbf{w}\|_1 + \frac{1-\beta}{2} \|\mathbf{w}\|^2 \right)$$

for some constant $\beta \in (0, 1)$

- Still gives a sparse \mathbf{w}
- Highly correlated variables will have similar values in \mathbf{w}

Outline

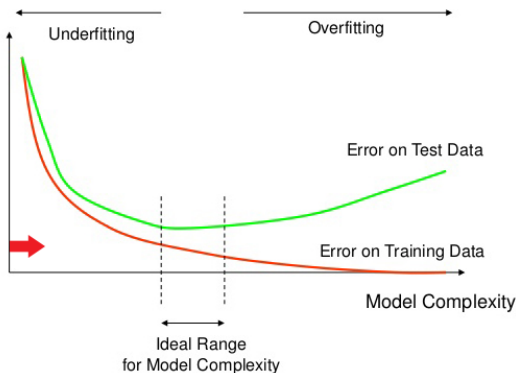
- 1 Learning Theory
- 2 Point Estimation: Bias and Variance
 - Consistency*
- 3 Decomposing Generalization Error
- 4 Regularization**
 - Weight Decay
 - Validation

Tuning Hyperparameters

- In ML, we call the constants that are fixed in a model the *hyperparameters*
 - Degree P in polynomial regression
 - Coefficient α of the weight decay term in the cost function of Ridge and LASSO, etc.
- Usually reflect some assumptions about the model
- Changing their values changes model complexity
 - And therefore generalization performance
- How to set appropriate values?
- Train a model many times with different hyperparameters, and choose the function with best generalizability
- Very time consuming, can we have heuristics to speed up the process?

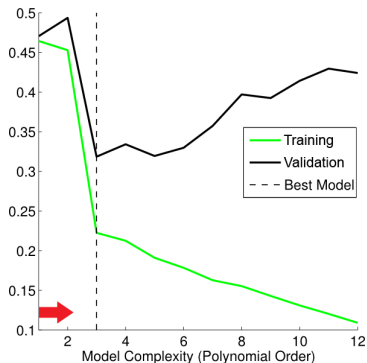
Structured Risk Minimization

- Consider again the Occam's razor
- Structured risk minimization**: start from the simplest model, gradually increase its complexity, and stop when overfitting



Validation Set

- Pitfall: we peep the testing set during the training process
 - The final function will overfit the testing set
 - Optimistic testing error
- Fix? Split a **validation set** from the training set and use it for hyperparameter selection



Reference I

- [1] Olivier Bousquet.
Concentration inequalities and empirical processes theory applied to the analysis of learning algorithms.
Ph.D. thesis, Ecole Polytechnique, Palaiseau, France, 2002.
- [2] Pascal Massart.
Some applications of concentration inequalities to statistics.
In *Annales de la Faculté des sciences de Toulouse: Mathématiques*, volume 9, pages 245–303, 2000.
- [3] Vladimir N Vapnik and A Ya Chervonenkis.
On the uniform convergence of relative frequencies of events to their probabilities.
In *Measures of Complexity*, pages 11–30. Springer, 2015.

Reference II

- [4] David H Wolpert.
The lack of a priori distinctions between learning algorithms.
Neural computation, 8(7):1341–1390, 1996.