

CS2109S Tutorial 5

Regularisation and Kernel Trick

(AY 25/26 Semester 2)

March 12, 2026

(Prepared by Benson)

Contents

Regularisation

Recap

Q1. L1-norm vs. L2-norm.

Q2. Normal Equation with L2 Regularization

Dual Representation and Kernel Trick

Recap: Dual Representation and Kernel Trick

Q3. Dual Representation of Logistic Regression

Q4. Kernel Trick

Bias and Variance

Recap

Extra. Bias and Variance

Bonus. Gaussian Kernel

Recap: Regularisation

- 📍 Why do we need regularisation?
 - A. To decrease model complexity.
 - B. To improve training time efficiency.
 - C. To avoid underfitting.
 - D. To avoid overfitting.

Recap: Regularisation

- 📍 How does regularisation address overfitting?
 - A. By penalizing weights for transformed features.
 - B. By reducing noise in the training data.
 - C. By penalizing large weights.
 - D. By reducing the number of transformed features.

Q1. L1-norm vs. L2-norm.

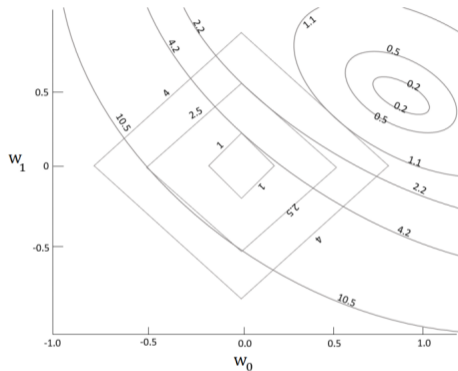
L2-norm (Ridge Regression):

$$J(\mathbf{w}) = \frac{1}{2m} \left[\sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^n \mathbf{w}_i^2$$

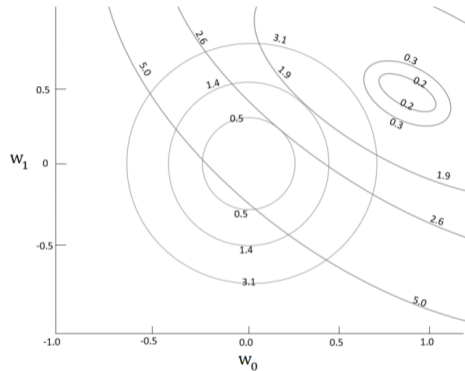
L1-norm (Lasso Regression):

$$J(\mathbf{w}) = \frac{1}{2m} \left[\sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^n |\mathbf{w}_i|$$

Q1. L1-norm vs. L2-norm.



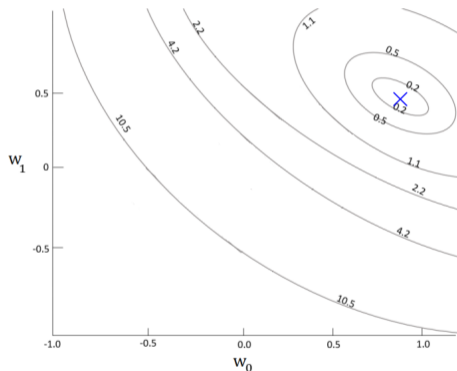
L1 regularisation



L2 regularisation

Q1. L1-norm vs. L2-norm.

(a) No regularisation.

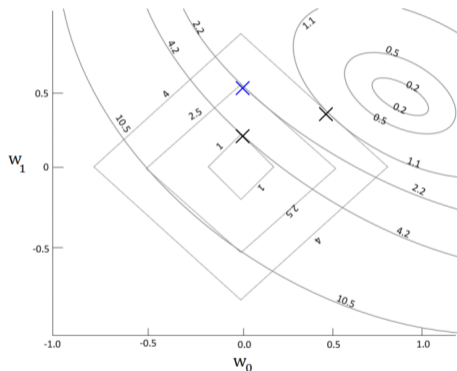


$$w_0 = 0.9, w_1 = 0.5$$

$$\text{Cost} \approx 0$$

Q1. L1-norm vs. L2-norm.

(b) L1 regularisation with $\lambda = 5$.



Total cost of the three points (from left to right):

▶ $4.2 + 1 = 5.2$

▶ $2.2 + 2.5 = 4.7$ 🏆

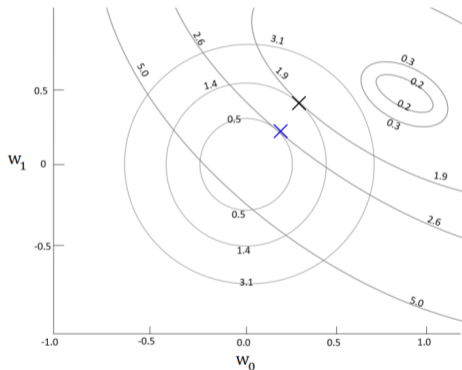
▶ $1.1 + 4 = 5.1$

$w_0 = 0.0, w_1 = 0.5$

Cost = 4.7

Q1. L1-norm vs. L2-norm.

(c) L2 regularisation with $\lambda = 5$.



Total cost of the two points (from left to right):

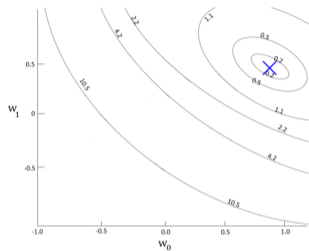
▶ $0.5 + 2.6 = 3.1$ 🏆

▶ $1.9 + 1.4 = 3.3$

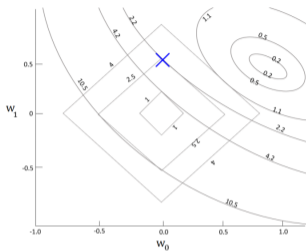
$w_0 = 0.2, w_1 = 0.25$

Cost = 3.1

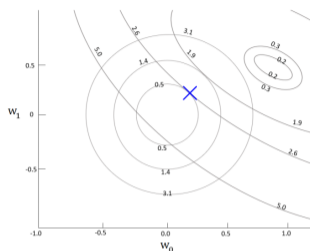
Q1. L1-norm vs. L2-norm.



No regularisation
 $w_0 = 0.9, w_1 = 0.5$
Cost ≈ 0



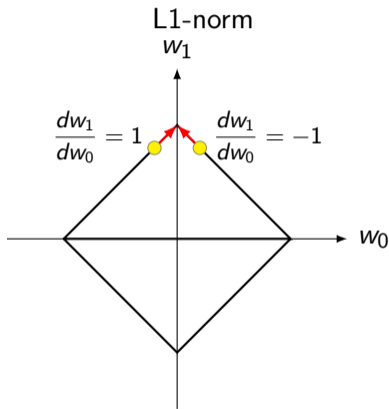
L1 regularisation
 $w_0 = 0.0, w_1 = 0.5$
Cost = 4.7



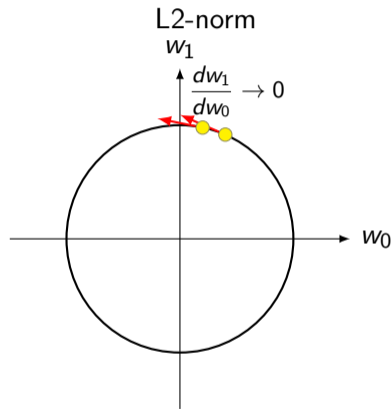
L2 regularisation
 $w_0 = 0.2, w_1 = 0.25$
Cost = 3.1

- ▶ L2 heavily penalizes larger parameters, preferring **all** smaller values.
- ▶ L1 may set values of certain parameters to 0 (why?).

Q1. L1-norm vs. L2-norm.



If w_1 is more important than w_0 , pushing towards w_1 is free lunch (“one for one”).

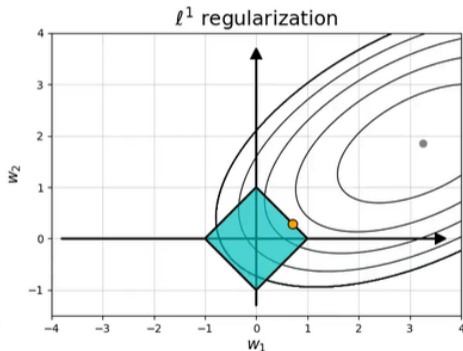
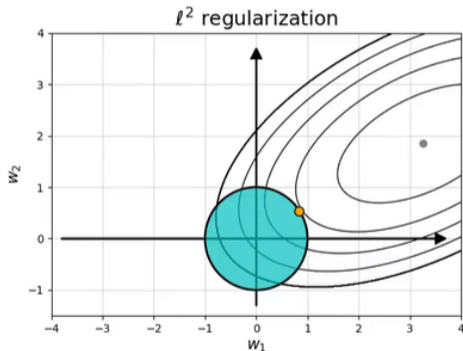


As $w_0 \rightarrow 0$, pushing towards w_1 has little gain (\therefore push will be “less aggressive”).

Q1. L1-norm vs. L2-norm.

Animation (See HTML slides):

ℓ^1 induces sparse solutions for least squares



by @itayevron

Q1. L1-norm vs. L2-norm.

L2-norm (Ridge Regression):

$$J(\mathbf{w}) = \frac{1}{2m} \left[\sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^n \mathbf{w}_i^2$$

- ▶ L2 heavily penalizes larger parameters, preferring **all** smaller values.

L1-norm (Lasso Regression):

$$J(\mathbf{w}) = \frac{1}{2m} \left[\sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \right] + \lambda \sum_{i=1}^n |\mathbf{w}_i|$$

- ▶ L1 may set values of certain parameters to 0 → effectively “feature selection” (select the more important features, and zero out the rest).

Q2. Normal Equation with L2 Regularization

- (a) Explain why the unnormalized normal equation cannot be used:

$$w = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ The two features x_1 and x_2 are linearly dependent. So $\mathbf{X}^T \mathbf{X}$ is not invertible.

x_1	x_2	y
1	2	3
2	4	5

Q2. Normal Equation with L2 Regularization

- (b) Suppose we add regularisation with parameter $\lambda = 1$.
Compute the weight vector \mathbf{w} .

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Solution.

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 4 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

x_1	x_2	y
1	2	3
2	4	5

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \begin{bmatrix} \frac{13}{33} & \frac{5}{11} & \frac{10}{11} \end{bmatrix}^T$$

$$\therefore \hat{y} = \frac{13}{33} + \frac{5}{11}x_1 + \frac{10}{11}x_2.$$

Note. $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is always invertible for all $\lambda > 0$.

Recap: Dual Representation and Kernel Trick

Primal Form of Linear Regression:

$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$$

Dual Form of Linear Regression:

$$h_{\alpha}(\mathbf{x}) = \left(\sum_{j=1}^n \alpha_j \mathbf{x}^{(j)} \right) \cdot \mathbf{x} = \sum_{j=1}^n \alpha_j (\mathbf{x}^{(j)} \cdot \mathbf{x})$$

What happens if we apply the **feature transformation** $\phi(\mathbf{x}) = [x \ x \ x \ x]^{\top}$?

Recap: Dual Representation and Kernel Trick

(Demo: See HTML slides)

Recap: Dual Representation and Kernel Trick

Important Property of a Kernel Function

$$K(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u}) \cdot \phi(\mathbf{v}) \text{ for some function } \phi$$

If we have a model that **ONLY** relies on the dot products of $\mathbf{x}^{(i)}$:

$$h_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j (\mathbf{x}^{(j)} \cdot \mathbf{x})$$

Applying the kernel function \equiv applying the transformed features!

$$h_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j (\phi(\mathbf{x}^{(j)}) \cdot \phi(\mathbf{x})) = \sum_{j=1}^n \alpha_j (k(\mathbf{x}^{(j)}, \mathbf{x}))$$

Q3. Dual Representation of Logistic Regression

Primal Form of Logistic Regression:

$$h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}$$

(a) Given $\mathbf{w} = \sum_{i=1}^n \alpha_i \mathbf{x}^{(i)}$, derive the dual representation of $h_{\mathbf{w}}(\mathbf{x})$.

$$\begin{aligned} h_{\alpha}(\mathbf{x}) &= \sigma \left(\left(\sum_{i=1}^n \alpha_i \mathbf{x}^{(i)} \right)^\top \mathbf{x} \right) \\ &= \sigma \left(\sum_{i=1}^n \alpha_i (\mathbf{x}^{(i)} \cdot \mathbf{x}) \right) \end{aligned}$$

(b) Write down the expression for the BCE loss $J^{BCE}(\alpha)$ in dual form.

$$J^{BCE}(\alpha) = \frac{1}{n} \sum_{i=1}^n BCE \left(y^{(i)}, \sigma \left(\sum_{i=1}^n \alpha_i (\mathbf{x}^{(i)} \cdot \mathbf{x}) \right) \right)$$

We do not need to explicitly compute the weight vector \mathbf{w} !

Q3. Dual Representation of Logistic Regression

Primal Form of Logistic Regression:

$$h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

Dual Form of Logistic Regression:

$$h_{\alpha}(\mathbf{x}) = \sigma \left(\sum_{i=1}^n \alpha_i (\mathbf{x}^{(i)} \cdot \mathbf{x}) \right)$$

- (c) How many parameters does the primal/dual logistic regression model have?
- ▶ Primal: d parameters.
 - ▶ Dual: n parameters.
- (d) When is primal/dual representation more efficient?
- ▶ Primal: When d is much smaller than n .
 - ▶ Dual: When d is large. Especially when the kernel trick could be used.

Q4. Kernel Trick

- (a) The linear kernel has a mapping $\phi(x) = x$. Compute $k(\mathbf{u}, \mathbf{v})$ for $\mathbf{u} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$ and $\mathbf{v} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$.

$$k(\mathbf{u}, \mathbf{v}) = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -1 \end{bmatrix} = -1$$

- (b) Consider the polynomial kernel of degree 2: $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^\top \mathbf{v} + 1)^2$. Compute $k(\mathbf{u}, \mathbf{v})$ for $\mathbf{u} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ and $\mathbf{v} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$.

$$k(\mathbf{u}, \mathbf{v}) = \left(\begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} + 1 \right)^2 = (5 + 1)^2 = 36$$

Q4. Kernel Trick

- (c) Consider the polynomial kernel of degree 2: $k(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^\top \mathbf{v} + 1)^2$. Find its explicit feature mapping.

$$\begin{aligned}k(\mathbf{u}, \mathbf{v}) &= (\mathbf{u}^\top \mathbf{v} + 1)^2 \\&= (u_1 v_1 + u_2 v_2 + 1)^2 \\&= u_1^2 v_1^2 + u_2^2 v_2^2 + 2u_1 u_2 v_1 v_2 + 2u_1 v_1 + 2u_2 v_2 + 1 \\&= \begin{bmatrix} u_1^2 \\ u_2^2 \\ \sqrt{2}u_1 u_2 \\ \sqrt{2}u_1 \\ \sqrt{2}u_2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} v_1^2 \\ v_2^2 \\ \sqrt{2}v_1 v_2 \\ \sqrt{2}v_1 \\ \sqrt{2}v_2 \\ 1 \end{bmatrix}\end{aligned}$$

Hence $\phi(\mathbf{x}) = [x_1^2 \quad x_2^2 \quad \sqrt{2}x_1 x_2 \quad \sqrt{2}x_1 \quad \sqrt{2}x_2 \quad 1]^\top$.

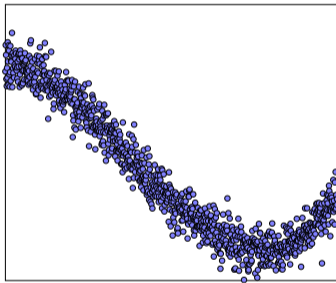
Q4. Kernel Trick

(d) What are the advantages of kernel trick?

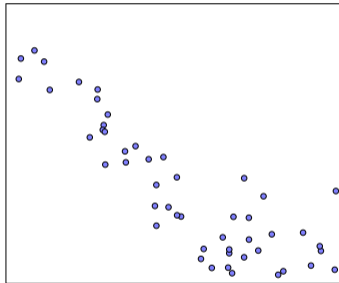
- ▶ We never explicitly compute $\phi(x)$, saving time and memory!
- ▶ We can change the feature transformation just by swapping $\phi(x)$.
- ▶ For infinite dimensional kernels (RBF / Gaussian kernel), explicit computation is impossible.

Bias and Variance

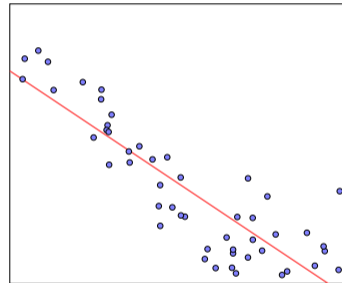
Population



Sample



Trained Model

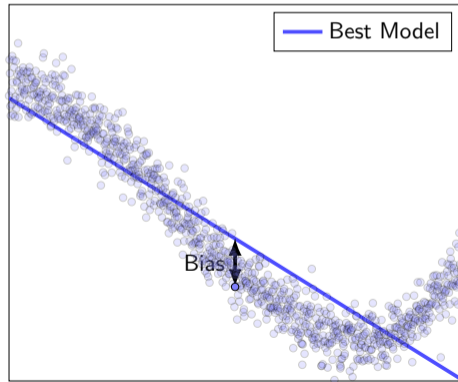
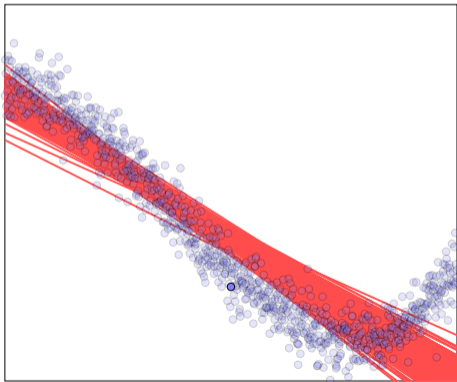


Hypothesis class:

$$h_w(x) = w_0x + b$$

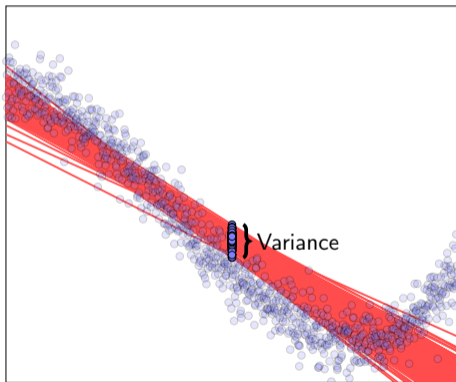
Bias and Variance

Bias measures how much the **expected predicted value** differs from the *actual value*.
Intuition: The error of the **best model** when you are given **infinite amount of training data**.



Bias and Variance

Variance measures the consistency of predictions due to sampling.



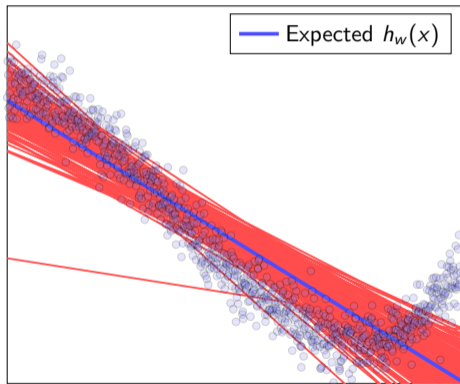
Bias and Variance

Poll:

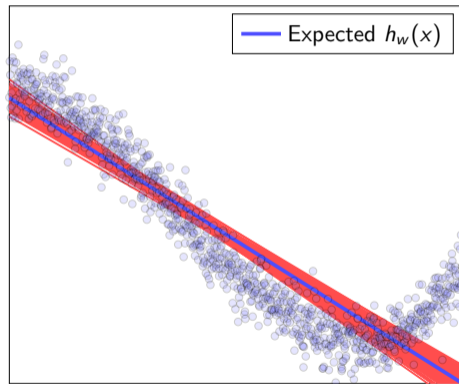
1. If we increase the amount of training samples:
The bias remains unchanged and the variance decreases.
2. If we increase the degree of the polynomial in our hypothesis class:
The bias decreases and the variance increases.

Bias and Variance

Less samples (Higher variance)

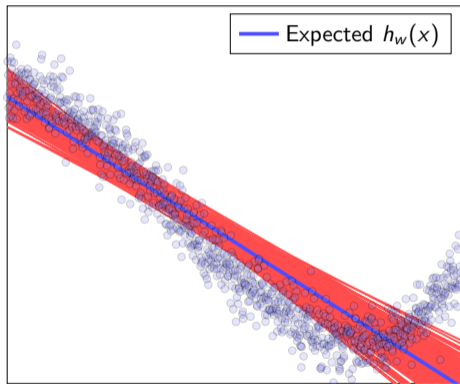


More samples (Lower variance)

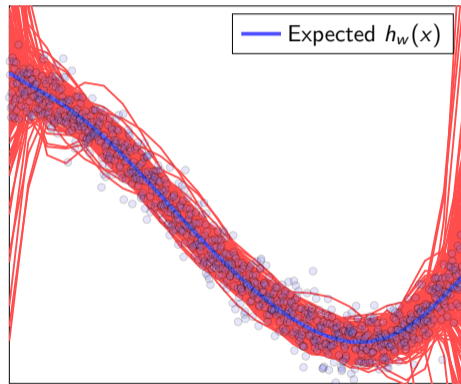


Bias and Variance

Lower degree (High bias; Low variance)

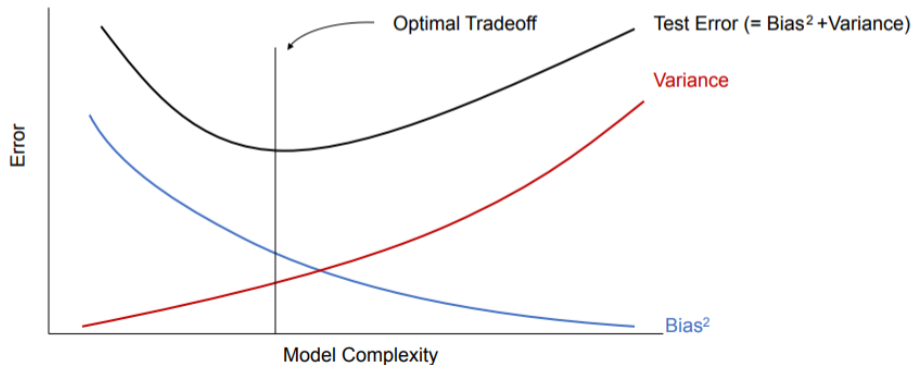


Higher degree (Low bias; High variance)



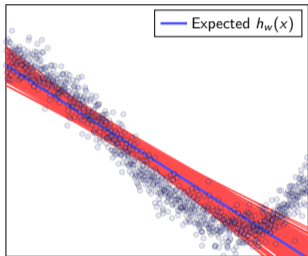
Bias-Variance Tradeoff

$$\text{Err}(x) = \underbrace{\left(y - \mathbb{E}[h_w(x)] \right)^2}_{\text{bias}} + \underbrace{\mathbb{E} \left[\left(\mathbb{E}[h_w(x)] - h_w(x) \right)^2 \right]}_{\text{variance}}$$

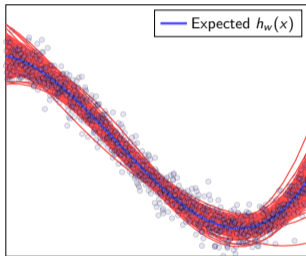


Bias-Variance Tradeoff

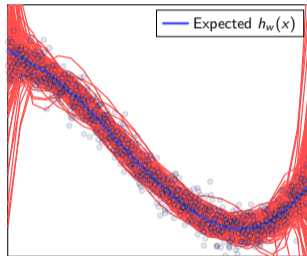
Underfitting



Just Right



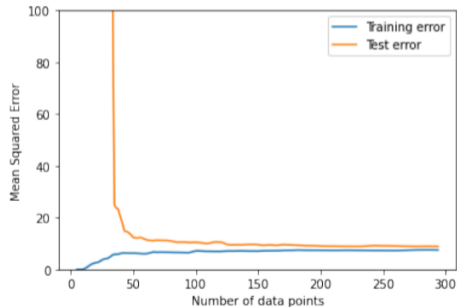
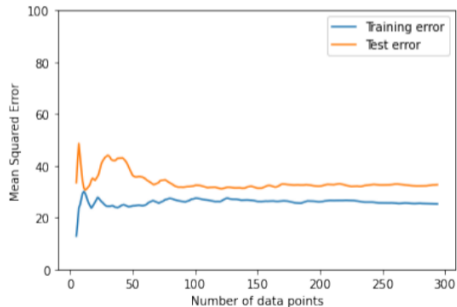
Overfitting



Extra. Bias and Variance

The model hypotheses are as below:

1. $H_w(x) = w_0 + w_1x$
2. $H_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_{10}x^{10}$



Extra. Bias and Variance

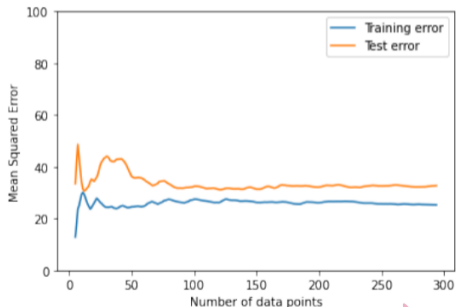
Main Idea:

- ▶ Both bias and variance contribute to the mean squared error.
- ▶ Bias does not decrease with the number of training samples. Bias causes error in **both** the training set and the test set.
- ▶ Variance decreases with the number of training samples. Variance causes a difference between the error in the training set and that in the test set.

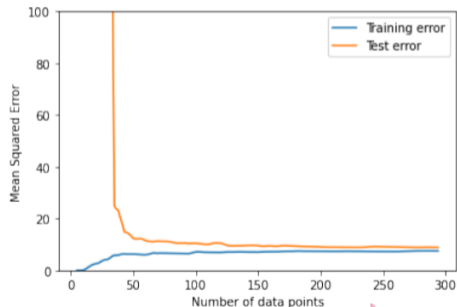
Q5. Bias and Variance

The model hypotheses are as below:

1. $H_w(x) = w_0 + w_1x$ **High bias**
2. $H_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_{10}x^{10}$ **High variance**



High bias



High variance

Bonus. Gaussian Kernel

Prove that the Gaussian Kernel, $K(\mathbf{u}, \mathbf{v}) = e^{-\frac{\|\mathbf{u}-\mathbf{v}\|^2}{2\sigma^2}}$, has infinite dimensional features.
(We assume $\sigma^2 = 1$ for simplicity.)

$$\begin{aligned}
 K(\mathbf{u}, \mathbf{v}) &= e^{-\frac{\|\mathbf{u}-\mathbf{v}\|^2}{2}} \\
 &= e^{-\frac{\mathbf{u}\cdot\mathbf{u}-2\mathbf{u}\cdot\mathbf{v}+\mathbf{v}\cdot\mathbf{v}}{2}} \\
 &= e^{-\frac{\mathbf{u}\cdot\mathbf{u}}{2}} e^{\mathbf{u}\cdot\mathbf{v}} e^{-\frac{\mathbf{v}\cdot\mathbf{v}}{2}} \\
 &= e^{-\frac{\mathbf{u}\cdot\mathbf{u}}{2}} \left(1 + \mathbf{u}\cdot\mathbf{v} + \frac{1}{2!}(\mathbf{u}\cdot\mathbf{v})^2 + \frac{1}{3!}(\mathbf{u}\cdot\mathbf{v})^3 + \frac{1}{4!}(\mathbf{u}\cdot\mathbf{v})^4 + \dots \right) e^{-\frac{\mathbf{v}\cdot\mathbf{v}}{2}} \\
 &= e^{-\frac{\mathbf{u}\cdot\mathbf{u}}{2}} (\phi_0(\mathbf{u}) \cdot \phi_0(\mathbf{v})) e^{-\frac{\mathbf{v}\cdot\mathbf{v}}{2}} \\
 &= \left(e^{-\frac{\mathbf{u}\cdot\mathbf{u}}{2}} \phi_0(\mathbf{u}) \right) \cdot \left(e^{-\frac{\mathbf{v}\cdot\mathbf{v}}{2}} \phi_0(\mathbf{v}) \right)
 \end{aligned}$$

$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$

Sum of polynomial kernels (in lecture)
 \Rightarrow Also a kernel (by bonus Q1)

scalars (multiplied to the transformed features)