

Bias / Variance Analysis, Kernel Methods

Professor Ameet Talwalkar

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Bias/Variance Analysis
- 4 Kernel methods

Announcements

- HW3 due now
- HW4 will online by Wednesday (due on 2/27)
- Midterm is on Wednesday

Midterm

- In-class from 10am - 11:50am
- Completely closed-book (no notes allowed)
- 6 short answer questions and 3 long questions
 - ▶ Short questions should take 5 minutes on average
 - ▶ Long questions should take 15 minutes each
- Covers all material through (and including) last Wednesday's class (i.e., today's review material)
 - ▶ Goal is to test conceptual understanding of the course material
 - ▶ Suggestion: carefully review lecture notes and problem sets

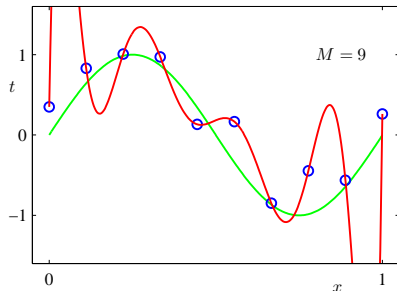
Outline

- 1 Administration
- 2 Review of last lecture
 - Basic ideas to combat overfitting
 - Ridge Regression and MAP estimate
- 3 Bias/Variance Analysis
- 4 Kernel methods

Overfitting

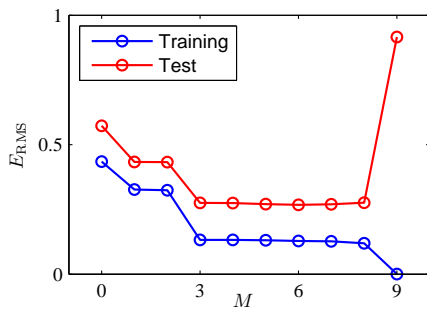
Example with regression, using polynomial basis functions

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$

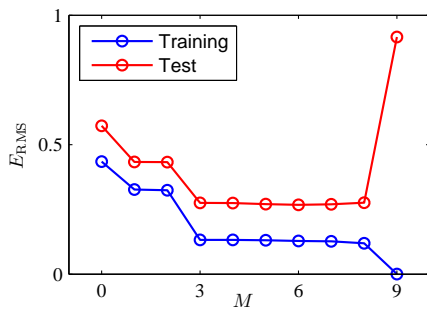


Can improve training accuracy at the expense of test accuracy

Visualizing overfitting



Visualizing overfitting

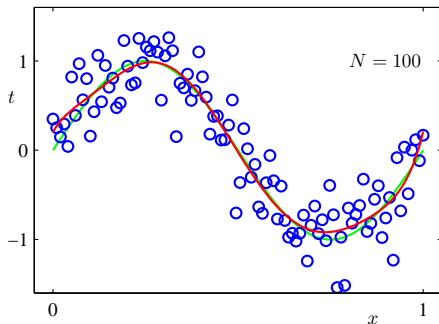


- As model becomes more complex, training error keeps improving while test error first improves then deteriorates

How to prevent overfitting?

How to prevent overfitting?

Use more training data



Regularization: adding a term to the objective function

$$\lambda \|\mathbf{w}\|_2^2$$

that favors a small parameter vector \mathbf{w} .

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}, \mathbf{w})$$

Probabilistic interpretation of Ridge Regression

- Ridge Regression model: $Y = \mathbf{w}^\top \mathbf{X} + \eta$
 - ▶ $Y \sim N(\mathbf{w}^\top \mathbf{X}, \sigma_0^2)$ is a Gaussian random variable (as before)
 - ▶ $w_d \sim N(0, \sigma^2)$ are i.i.d. Gaussian random variables (unlike before)
 - ▶ Note that all w_d share the same variance σ^2
- \mathbf{w} is a random variable (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

- Maximum a posterior (MAP) estimate:

$$\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathcal{D}) = \arg \max_{\mathbf{w}} p(\mathcal{D}, \mathbf{w})$$

- MAP reduces to MLE if we assume uniform prior for $p(\mathbf{w})$

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Joint log likelihood

Plugging in Gaussian PDF, we get:

$$\begin{aligned} \log p(\mathcal{D}, \mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

Estimating \mathbf{w}

- Let X_1, \dots, X_N be IID with $y|\mathbf{w}, \mathbf{x} \sim N(\mathbf{w}^\top \mathbf{x}, \sigma_0^2)$
- Let w_d be IID with $w_d \sim N(0, \sigma^2)$

Joint likelihood of data and parameters (given σ_0, σ)

$$p(\mathcal{D}, \mathbf{w}) = p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) = \prod_n p(y_n|\mathbf{x}_n, \mathbf{w}) \prod_d p(w_d)$$

Joint log likelihood

Plugging in Gaussian PDF, we get:

$$\begin{aligned} \log p(\mathcal{D}, \mathbf{w}) &= \sum_n \log p(y_n|\mathbf{x}_n, \mathbf{w}) + \sum_d \log p(w_d) \\ &= -\frac{\sum_n (\mathbf{w}^\top \mathbf{x}_n - y_n)^2}{2\sigma_0^2} - \sum_d \frac{1}{2\sigma^2} w_d^2 + \text{const} \end{aligned}$$

MAP estimate: $\mathbf{w}^{\text{MAP}} = \arg \max_{\mathbf{w}} \log p(\mathcal{D}, \mathbf{w})$

- As with LMS, set gradient equal to zero and solve (for \mathbf{w})

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

$$\mathcal{E}(\mathbf{w}) = \sum_n (\mathbf{w}^T \mathbf{x}_n - y_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$ is used to denote σ_0^2/σ^2 .

- What happens as $\lambda \rightarrow +\infty$?

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

$$\mathcal{E}(\mathbf{w}) = \sum_n (\mathbf{w}^T \mathbf{x}_n - y_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$ is used to denote σ_0^2/σ^2 .

- What happens as $\lambda \rightarrow +\infty$?
- What happens as $\lambda \rightarrow 0$?

Maximum a posterior (MAP) estimate

Regularized linear regression: a new error to minimize

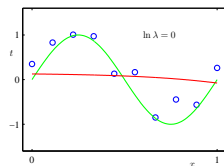
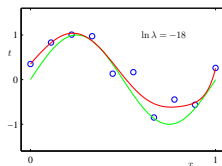
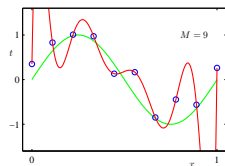
$$\mathcal{E}(\mathbf{w}) = \sum_n (\mathbf{w}^T \mathbf{x}_n - y_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$ is used to denote σ_0^2/σ^2 .

- What happens as $\lambda \rightarrow +\infty$?
- What happens as $\lambda \rightarrow 0$?
- What happens when $\lambda < 0$?

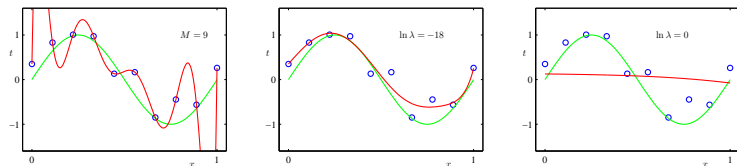
Overfitting in terms of λ

Overfitting is reduced as we increase the regularizer

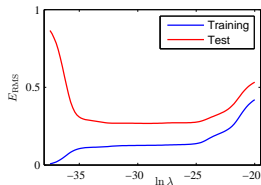


Overfitting in terms of λ

Overfitting is reduced as we increase the regularizer



λ vs. **residual error** shows the difference of the model performance on training and testing dataset



Outline

- 1 Administration
- 2 Review of last lecture
- 3 Bias/Variance Analysis**
- 4 Kernel methods

Basic and important machine learning concepts

Supervised learning

We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(h(\mathbf{x}), y)$$

Basic and important machine learning concepts

Supervised learning

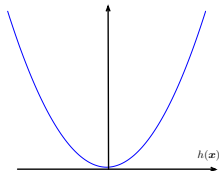
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

$$\ell(h(\mathbf{x}), y)$$

Example: quadratic loss function for regression when y is continuous

$$\ell(h(\mathbf{x}), y) = [h(\mathbf{x}) - y]^2$$

Ex: when $y = 0$

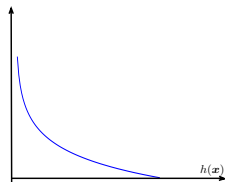


Other types of loss functions

For classification: cross-entropy loss (also called *logistic* loss)

$$\ell(h(\mathbf{x}), y) = -y \log h(\mathbf{x}) - (1-y) \log[1-h(\mathbf{x})]$$

Ex: when $y = 1$



Measure how good our predictor is

Risk: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h(\mathbf{x}), y) = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

Measure how good our predictor is

Risk: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h(\mathbf{x}), y) = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

However, we cannot compute $R[h(\mathbf{x})]$, so we use *empirical risk*, given a training dataset \mathcal{D}

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Measure how good our predictor is

Risk: assume we know the true distribution of data $p(\mathbf{x}, y)$, the *risk* is

$$R[h(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h(\mathbf{x}), y) = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} dy$$

However, we cannot compute $R[h(\mathbf{x})]$, so we use *empirical risk*, given a training dataset \mathcal{D}

$$R^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(h(\mathbf{x}_n), y_n)$$

Intuitively, as $N \rightarrow +\infty$,

$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow R[h(\mathbf{x})]$$

How does this relate to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

How does this relate to what we have learned?

So far, we have been doing empirical risk minimization (ERM)

- For linear regression, $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$, and we use squared loss
- For logistic regression, $h(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$, and we use cross-entropy loss

ERM might be problematic

- If $h(\mathbf{x})$ is complicated enough,

$$R^{\text{EMP}}[h(\mathbf{x})] \rightarrow 0$$

- But then $h(\mathbf{x})$ is unlikely to do well in predicting things out of the training dataset \mathcal{D} (*poor generalization* or *overfitting*)
- We'll explore why regularization might work from the context of the bias-variance tradeoff, focusing on regression / squared loss

Bias/variance tradeoff (Looking ahead)

Error decomposes into 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

We will prove this result, and interpret what it means...

Bias/variance tradeoff for regression

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data

Bias/variance tradeoff for regression

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data
- $h_{\mathcal{D}}(\mathbf{x})$: our prediction function
 - ▶ subscript \mathcal{D} indicates that function is learned on a specific training set

Bias/variance tradeoff for regression

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data
- $h_{\mathcal{D}}(\mathbf{x})$: our prediction function
 - ▶ subscript \mathcal{D} indicates that function is learned on a specific training set
- $\ell(h(\mathbf{x}), y)$: squared loss function for regression

$$\ell(h_{\mathcal{D}}(\mathbf{x}), y) = [h_{\mathcal{D}}(\mathbf{x}) - y]^2$$

Bias/variance tradeoff for regression

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data
- $h_{\mathcal{D}}(\mathbf{x})$: our prediction function
 - ▶ subscript \mathcal{D} indicates that function is learned on a specific training set
- $\ell(h(\mathbf{x}), y)$: squared loss function for regression

$$\ell(h_{\mathcal{D}}(\mathbf{x}), y) = [h_{\mathcal{D}}(\mathbf{x}) - y]^2$$

- Unknown joint distribution $p(\mathbf{x}, y)$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

But, \mathcal{D} is a random sample from the following joint distribution

$$\mathcal{D} \sim P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n)$$

The effect of finite training samples

Risk of prediction function, $h_{\mathcal{D}}(\mathbf{x})$

$$R[h_{\mathcal{D}}(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} \ell(h_{\mathcal{D}}(\mathbf{x}), y) = \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

But, \mathcal{D} is a random sample from the following joint distribution

$$\mathcal{D} \sim P(\mathcal{D}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n)$$

So, $h_{\mathcal{D}}(\mathbf{x})$ and $R[h_{\mathcal{D}}(\mathbf{x})]$ are also random w.r.t. $P(\mathcal{D})$

How can we disentangle the impact of the random sample \mathcal{D} when assessing the quality of $h_{\mathcal{D}}(\cdot)$?

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to \mathcal{D} is marginalized out.

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to \mathcal{D} is marginalized out.

Averaged prediction

$$\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) = \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D}$$

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Namely, the randomness with respect to \mathcal{D} is marginalized out.

Averaged prediction

$$\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) = \int_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) P(\mathcal{D}) d\mathcal{D}$$

Namely, if we have seen many training datasets, we predict with the average of our trained models learned on each training dataset.

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y \left[\left(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) \right) \right. \\ &\quad \left. + \left(\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y \right) \right]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y \left[\left(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) \right) \right. \\ &\quad \left. + \left(\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y \right) \right]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE}}\end{aligned}$$

Interpreting Averaged risk

We can add and subtract averaged prediction from averaged risk

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y \left[\left(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) \right) \right. \\ &\quad \left. + \left(\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y \right) \right]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \underbrace{\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}}_{\text{VARIANCE}} \\ &\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

Where does the cross-term go?

It is zero

$$\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})][\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Where does the cross-term go?

It is zero

$$\begin{aligned} & \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})][\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathbf{x}} \int_y \left\{ \int_{\mathcal{D}} [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})] P(\mathcal{D}) d\mathcal{D} \right\} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy \end{aligned}$$

Where does the cross-term go?

It is zero

$$\begin{aligned} & \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})][\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &= \int_{\mathbf{x}} \int_y \left\{ \int_{\mathcal{D}} [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})] P(\mathcal{D}) d\mathcal{D} \right\} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y] p(\mathbf{x}, y) d\mathbf{x} dy \\ &= 0 \leftarrow \text{(the integral within the braces vanishes, by definition)} \end{aligned}$$

Analyzing the variance

How can we reduce the variance?

$$\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

Analyzing the variance

How can we reduce the variance?

$$\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

- Use a lot of data (ie, increase the size of \mathcal{D})
- Use a simple $h(\cdot)$ so that $h_{\mathcal{D}}(\mathbf{x})$ does not vary much across different training datasets.

Analyzing the variance

How can we reduce the variance?

$$\int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}$$

- Use a lot of data (ie, increase the size of \mathcal{D})
- Use a simple $h(\cdot)$ so that $h_{\mathcal{D}}(\mathbf{x})$ does not vary much across different training datasets.

Ex: $h(\mathbf{x}) = \text{const}$

The remaining item

$$\begin{aligned}\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

The remaining item

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

The integrand has no dependency on \mathcal{D} anymore and simplifies to

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

The remaining item

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}R[h_{\mathcal{D}}(\mathbf{x})] &= \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x})]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D} \\ &\quad + \int_{\mathcal{D}} \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy P(\mathcal{D}) d\mathcal{D}\end{aligned}$$

The integrand has no dependency on \mathcal{D} anymore and simplifies to

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}}h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

We will apply a similar trick, by using the averaged target y

$$\mathbb{E}_y[y] = \int_y yp(y|\mathbf{x})dy$$

Bias and noise

Decompose again

$$\begin{aligned} & \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \int_{\mathbf{x}} \int_y \left[\left(\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y] \right) + \left(\mathbb{E}_y[y] - y \right) \right]^2 p(\mathbf{x}, y) d\mathbf{x} dy \end{aligned}$$

Bias and noise

Decompose again

$$\begin{aligned} & \int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \int_{\mathbf{x}} \int_y \left[\left(\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y] \right) + \left(\mathbb{E}_y[y] - y \right) \right]^2 p(\mathbf{x}, y) d\mathbf{x} dy \\ &= \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{BIAS}^2} \\ & \quad + \underbrace{\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy}_{\text{NOISE}} \end{aligned}$$

Where is the cross-term?

Left as a take-home exercise

Analyzing the noise

How can we reduce noise

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left(\int_y [\mathbb{E}_y[y] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

Analyzing the noise

How can we reduce noise

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left(\int_y [\mathbb{E}_y[y] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

There is *nothing* we can do. This quantity depends on $p(\mathbf{x}, y)$ only; choosing $h(\cdot)$ or the training dataset \mathcal{D} will not affect it.

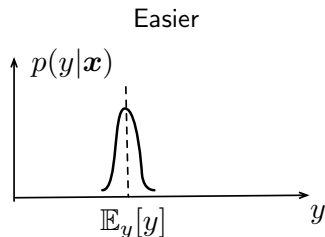
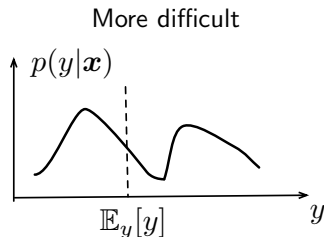
Analyzing the noise

How can we reduce noise

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left(\int_y [\mathbb{E}_y[y] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

There is *nothing* we can do. This quantity depends on $p(\mathbf{x}, y)$ only; choosing $h(\cdot)$ or the training dataset \mathcal{D} will not affect it.

Note that the integral inside the parentheses is the *variance* (noise) of the distribution $p(y|\mathbf{x})$ at the given \mathbf{x} .



Analyzing the bias term

How can we reduce bias?

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}) d\mathbf{x}$$

Analyzing the bias term

How can we reduce bias?

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}) d\mathbf{x}$$

By using richer / more complex models to better approximate $\mathbb{E}_y[y]$

Analyzing the bias term

How can we reduce bias?

$$\int_{\mathbf{x}} \int_y [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y]]^2 p(\mathbf{x}) d\mathbf{x}$$

By using richer / more complex models to better approximate $\mathbb{E}_y[y]$

However, this increased complexity will increase the VARIANCE term (as we can potentially overfit)

Bias/variance tradeoff

Error decomposes into 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

where the first and the second term are inherently in conflict in terms of choosing what kind of $h(\mathbf{x})$ we should use (unless we have an infinite amount of data)

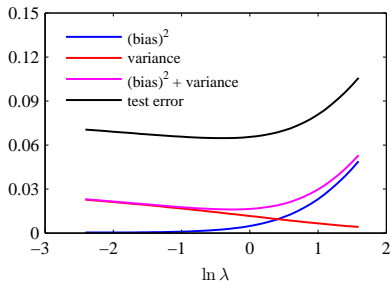
Bias/variance tradeoff

Error decomposes into 3 terms

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})] = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}$$

where the first and the second term are inherently in conflict in terms of choosing what kind of $h(\mathbf{x})$ we should use (unless we have an infinite amount of data)

If we can compute all terms analytically, they will look like this



Outline

- 1 Administration
- 2 Review of last lecture
- 3 Bias/Variance Analysis
- 4 Kernel methods**
 - Motivation
 - Kernel matrix and kernel functions
 - Kernelized machine learning methods

Motivation

How to choose nonlinear basis function for regression?

$$\mathbf{w}^T \phi(\mathbf{x})$$

- $\phi(\cdot)$ maps the original feature vector \mathbf{x} to a *new* M -dimensional feature vector
- We can sidestep the issue of choosing which $\phi(\cdot)$ to use by *equivalently* choosing a *kernel function*

Motivation

How to choose nonlinear basis function for regression?

$$\mathbf{w}^T \phi(\mathbf{x})$$

- $\phi(\cdot)$ maps the original feature vector \mathbf{x} to a *new* M -dimensional feature vector
- We can sidestep the issue of choosing which $\phi(\cdot)$ to use by *equivalently* choosing a *kernel function*

Regularized least squares

$$J(\mathbf{w}) = \frac{1}{2} \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Its solution \mathbf{w}^{MAP} is given by

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))(-\phi(\mathbf{x}_n)) + \lambda \mathbf{w} = 0$$

MAP Solution

The optimal parameter vector is a linear combination of features

$$\mathbf{w}^{\text{MAP}} = \sum_n \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)$$

MAP Solution

The optimal parameter vector is a linear combination of features

$$\mathbf{w}^{\text{MAP}} = \sum_n \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) = \sum_n \alpha_n \phi(\mathbf{x}_n) = \Phi^T \boldsymbol{\alpha}$$

- $\alpha_n = \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$
- Φ is the *design matrix* of *transformed* features
- Its transpose is made of column vectors and is given by

$$\Phi^T = (\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \ \cdots \ \phi(\mathbf{x}_N)) \in \mathbb{R}^{M \times N}$$

MAP Solution

The optimal parameter vector is a linear combination of features

$$\mathbf{w}^{\text{MAP}} = \sum_n \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n) = \sum_n \alpha_n \phi(\mathbf{x}_n) = \Phi^T \boldsymbol{\alpha}$$

- $\alpha_n = \frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$
- Φ is the *design matrix* of transformed features
- Its transpose is made of column vectors and is given by

$$\Phi^T = (\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \ \cdots \ \phi(\mathbf{x}_N)) \in \mathbb{R}^{M \times N}$$

Of course, we don't know what $\boldsymbol{\alpha}$ (the vector of α_n) corresponds to \mathbf{w}^{MAP} !

Dual formulation

Regularized least squares: $J(\mathbf{w}) = \frac{1}{2} \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$

Substitute $\mathbf{w}^{\text{MAP}} = \Phi^T \alpha$ into $J(\mathbf{w})$ to obtain a function of α :

$$J(\alpha) = \frac{1}{2} \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha - (\Phi \Phi^T \mathbf{y})^T \alpha + \frac{\lambda}{2} \alpha^T \Phi \Phi^T \alpha$$

Dual formulation

Regularized least squares: $J(\mathbf{w}) = \frac{1}{2} \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$

Substitute $\mathbf{w}^{\text{MAP}} = \Phi^T \alpha$ into $J(\mathbf{w})$ to obtain a function of α :

$$J(\alpha) = \frac{1}{2} \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha - (\Phi \Phi^T \mathbf{y})^T \alpha + \frac{\lambda}{2} \alpha^T \Phi \Phi^T \alpha$$

Before we show how $J(\alpha)$ is derived, we make an important observation

Dual formulation

Regularized least squares: $J(\mathbf{w}) = \frac{1}{2} \sum_n (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$

Substitute $\mathbf{w}^{\text{MAP}} = \Phi^T \alpha$ into $J(\mathbf{w})$ to obtain a function of α :

$$J(\alpha) = \frac{1}{2} \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha - (\Phi \Phi^T \mathbf{y})^T \alpha + \frac{\lambda}{2} \alpha^T \Phi \Phi^T \alpha$$

Before we show how $J(\alpha)$ is derived, we make an important observation

The *Gram matrix* or *kernel matrix* — $\Phi \Phi^T$ — appears multiple times

$$\mathbf{K} = \Phi \Phi^T$$

$$= \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_1)^T \phi(\mathbf{x}_N) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_2)^T \phi(\mathbf{x}_N) \\ \cdots & \cdots & \cdots & \cdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_1) & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_2) & \cdots & \phi(\mathbf{x}_N)^T \phi(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times N}$$

Properties of the matrix \mathbf{K}

- Symmetric

$$K_{mn} = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = K_{nm}$$

- Positive semidefinite: for any vector \mathbf{a}

$$\mathbf{a}^T \mathbf{K} \mathbf{a} = (\Phi^T \mathbf{a})^T (\Phi^T \mathbf{a}) \geq 0$$

Properties of the matrix \mathbf{K}

- Symmetric

$$K_{mn} = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = K_{nm}$$

- Positive semidefinite: for any vector \mathbf{a}

$$\mathbf{a}^T \mathbf{K} \mathbf{a} = (\Phi^T \mathbf{a})^T (\Phi^T \mathbf{a}) \geq 0$$

- Not the same as the second-moment (covariance) matrix $\mathbf{C} = \Phi^T \Phi$
 - ▶ \mathbf{C} has a size of $M \times M$ while \mathbf{K} is $N \times N$
 - ▶ When $N \leq M$, using \mathbf{K} is more computationally advantageous

The derivation of $J(\boldsymbol{\alpha})$

$$J(\mathbf{w}) = \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \end{aligned} \quad (\mathbf{y}, \boldsymbol{\Phi} \mathbf{w} \in \mathbb{R}^N)$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 && (\mathbf{y}, \boldsymbol{\Phi} \mathbf{w} \in \mathbb{R}^N) \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 && (\mathbf{w}^{\text{MAP}} = \boldsymbol{\Phi}^T \boldsymbol{\alpha}) \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 && (\mathbf{y}, \boldsymbol{\Phi} \mathbf{w} \in \mathbb{R}^N) \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 && (\mathbf{w}^{\text{MAP}} = \boldsymbol{\Phi}^T \boldsymbol{\alpha}) \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} && (\mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^T, \|\mathbf{v}\|_2^2 = \mathbf{v}^T \mathbf{v}) \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 && (\mathbf{y}, \boldsymbol{\Phi} \mathbf{w} \in \mathbb{R}^N) \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 && (\mathbf{w}^{\text{MAP}} = \boldsymbol{\Phi}^T \boldsymbol{\alpha}) \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} && (\mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^T, \|\mathbf{v}\|_2^2 = \mathbf{v}^T \mathbf{v}) \\ &\propto \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\alpha} - \mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \end{aligned}$$

The derivation of $J(\boldsymbol{\alpha})$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 && (\mathbf{y}, \boldsymbol{\Phi} \mathbf{w} \in \mathbb{R}^N) \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 && (\mathbf{w}^{\text{MAP}} = \boldsymbol{\Phi}^T \boldsymbol{\alpha}) \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} && (\mathbf{K} = \boldsymbol{\Phi} \boldsymbol{\Phi}^T, \|\mathbf{v}\|_2^2 = \mathbf{v}^T \mathbf{v}) \\ &\propto \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\alpha} - \mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \\ &= \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K}^2 \boldsymbol{\alpha} - (\mathbf{K} \mathbf{y})^T \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} && (\mathbf{K} \text{ is symmetric}) \\ &= J(\boldsymbol{\alpha}) \end{aligned}$$

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} \mathbf{y} + \lambda \mathbf{K} \alpha = 0$$

which leads to (assuming that \mathbf{K} is invertible)

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

- We only need to know \mathbf{K} in order to compute α !
- Solution doesn't involve $\phi(\cdot)$, but instead inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$
- This observation will give rise to the use of *kernel functions*

Optimal α

$$\frac{\partial J(\alpha)}{\partial \alpha} = \mathbf{K}^2 \alpha - \mathbf{K} \mathbf{y} + \lambda \mathbf{K} \alpha = 0$$

which leads to (assuming that \mathbf{K} is invertible)

$$\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

- We only need to know \mathbf{K} in order to compute α !
- Solution doesn't involve $\phi(\cdot)$, but instead inner products $\phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$
- This observation will give rise to the use of *kernel functions*

Note that the initial parameter vector does require knowledge of Φ

$$\mathbf{w}^{\text{MAP}} = \Phi^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Computing prediction needs only inner products too!

Since $\mathbf{w}^{\text{MAP}} = \Phi^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$, at test time we compute:

$$\mathbf{w}^T\phi(\mathbf{x}) = \mathbf{y}^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\Phi\phi(\mathbf{x})$$

Computing prediction needs only inner products too!

Since $\mathbf{w}^{\text{MAP}} = \Phi^T(\mathbf{K} + \lambda\mathbf{I})^{-1}\mathbf{y}$, at test time we compute:

$$\begin{aligned} \mathbf{w}^T \phi(\mathbf{x}) &= \mathbf{y}^T (\mathbf{K} + \lambda\mathbf{I})^{-1} \Phi \phi(\mathbf{x}) \\ &= \mathbf{y}^T (\mathbf{K} + \lambda\mathbf{I})^{-1} \begin{pmatrix} \phi(\mathbf{x}_1)^T \phi(\mathbf{x}) \\ \phi(\mathbf{x}_2)^T \phi(\mathbf{x}) \\ \vdots \\ \phi(\mathbf{x}_N)^T \phi(\mathbf{x}) \end{pmatrix} = \mathbf{y}^T (\mathbf{K} + \lambda\mathbf{I})^{-1} \mathbf{k}_x \end{aligned}$$

- We used the property that $(\mathbf{K} + \lambda\mathbf{I})^{-1}$ is symmetric (as \mathbf{K} is)
- \mathbf{k}_x is shorthand notation for the column vector of inner products between training set and test point
- To make a prediction we *only need* $\phi(\mathbf{x}_n)^T \phi(\mathbf{x})!$

E.g.: Dot product between degree two polynomial features

Consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

E.g.: Dot product between degree two polynomial features

Consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

This gives rise to an inner product in a special form,

$$\begin{aligned} \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) &= x_{m1}^2 x_{n1}^2 + 2x_{m1}x_{m2}x_{n1}x_{n2} + x_{m2}^2 x_{n2}^2 \\ &= (x_{m1}x_{n1} + x_{m2}x_{n2})^2 = (\mathbf{x}_m^T \mathbf{x}_n)^2 \end{aligned}$$

E.g.: Dot product between degree two polynomial features

Consider the following $\phi(\mathbf{x})$:

$$\phi : \mathbf{x} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

This gives rise to an inner product in a special form,

$$\begin{aligned} \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) &= x_{m1}^2 x_{n1}^2 + 2x_{m1}x_{m2}x_{n1}x_{n2} + x_{m2}^2 x_{n2}^2 \\ &= (x_{m1}x_{n1} + x_{m2}x_{n2})^2 = (\mathbf{x}_m^T \mathbf{x}_n)^2 \end{aligned}$$

Namely, inner product can be computed by a function $(\mathbf{x}_m^T \mathbf{x}_n)^2$ defined in terms of the original features, *without knowing $\phi(\cdot)$!*

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $c \geq 0$ and d is a positive integer.

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Shift-invariant kernel (only depends on difference between two inputs)
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)

Common kernel functions

Polynomial kernel function with degree of d

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n + c)^d$$

for $c \geq 0$ and d is a positive integer.

Gaussian kernel, RBF kernel, or Gaussian RBF kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

- Shift-invariant kernel (only depends on difference between two inputs)
- Corresponds to a feature space with *infinite* dimensions (but we can work directly with the original features)

These kernels have hyperparameters to be tuned: d , c , σ^2

Kernel functions

Definition: a (positive semidefinite) kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Kernel functions

Definition: a (positive semidefinite) kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Kernel functions

Definition: a (positive semidefinite) kernel function $k(\cdot, \cdot)$ is a bivariate function that satisfies the following properties. For any \mathbf{x}_m and \mathbf{x}_n ,

$$k(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_n, \mathbf{x}_m) \text{ and } k(\mathbf{x}_m, \mathbf{x}_n) = \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$$

for *some* function $\phi(\cdot)$.

Examples we have seen

$$k(\mathbf{x}_m, \mathbf{x}_n) = (\mathbf{x}_m^T \mathbf{x}_n)^2$$

$$k(\mathbf{x}_m, \mathbf{x}_n) = e^{-\|\mathbf{x}_m - \mathbf{x}_n\|_2^2 / 2\sigma^2}$$

Example that is not a kernel

$$k(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2$$

(we'll see why later)

Conditions for being a positive semidefinite kernel function

Mercer theorem (loosely), a bivariate function $k(\cdot, \cdot)$ is a positive semidefinite kernel function, if and only if, for *any* N and *any* $\mathbf{x}_1, \mathbf{x}_2, \dots$, and \mathbf{x}_N , the matrix

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

is positive semidefinite. We also refer $k(\cdot, \cdot)$ as a positive semidefinite kernel.

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

Why $\|\mathbf{x}_m - \mathbf{x}_n\|_2^2$ is not a positive semidefinite kernel?

Use the definition of positive semidefinite kernel function. We choose $N = 2$, and compute the matrix

$$\mathbf{K} = \begin{pmatrix} 0 & \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \\ \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 & 0 \end{pmatrix}$$

- PSD matrices have only non-negative eigenvalues
- This matrix has both *negative* and positive eigenvalues
- *Trace* of a matrix equals the sum of the diagonal elements, and also equals to the sum of the matrix's eigenvalues
 - ▶ In our case, the trace is zero

Why use kernel functions?

Can define kernel matrix without specifying $\phi(\cdot)$

$$\mathbf{K} = \Phi\Phi^T = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \vdots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

'Kernel trick'

- Many learning methods rely on training and test data only in the form of *inner products*, e.g., regularized least squares, nearest neighbors
- We can use a kernel function to introduce nonlinearity, i.e., *"kernelizing"* the methods
- We will show this "trick" by kernelizing the nearest neighbor classifier
- We will see this again when we talk about support vector machines

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Can replace dot products by a kernel function $k(\cdot, \cdot)$:

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n)$$

Kernelized nearest neighbors classifier (NNC)

Fundamental quantity is (squared) distance between two data points

$$d(\mathbf{x}_m, \mathbf{x}_n) = \|\mathbf{x}_m - \mathbf{x}_n\|_2^2 = \mathbf{x}_m^T \mathbf{x}_m + \mathbf{x}_n^T \mathbf{x}_n - 2\mathbf{x}_m^T \mathbf{x}_n$$

Can replace dot products by a kernel function $k(\cdot, \cdot)$:

$$d^{\text{KERNEL}}(\mathbf{x}_m, \mathbf{x}_n) = k(\mathbf{x}_m, \mathbf{x}_m) + k(\mathbf{x}_n, \mathbf{x}_n) - 2k(\mathbf{x}_m, \mathbf{x}_n) = d(\phi(\mathbf{x}_m), \phi(\mathbf{x}_n))$$

- d^{KERNEL} equivalent to distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$ where $\phi(\cdot)$ is the nonlinear mapping implied by the kernel function
- The nearest neighbor of a test point \mathbf{x} is found via

$$\arg \min_n d^{\text{KERNEL}}(\mathbf{x}, \mathbf{x}_n)$$

There are infinite numbers of kernels to use!

Rules of composing kernels (this is just a partial list)

- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $ck(\mathbf{x}_m, \mathbf{x}_n)$ is also if $c > 0$.
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $\alpha k_1(\mathbf{x}_m, \mathbf{x}_n) + \beta k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also if $\alpha, \beta \geq 0$
- if both $k_1(\mathbf{x}_m, \mathbf{x}_n)$ and $k_2(\mathbf{x}_m, \mathbf{x}_n)$ are kernels, then $k_1(\mathbf{x}_m, \mathbf{x}_n)k_2(\mathbf{x}_m, \mathbf{x}_n)$ are also.
- if $k(\mathbf{x}_m, \mathbf{x}_n)$ is a kernel, then $e^{k(\mathbf{x}_m, \mathbf{x}_n)}$ is also.
- ...

In practice, choosing an appropriate kernel is an “art”

People typically start with polynomial and Gaussian RBF kernels or incorporate domain knowledge.