

Bias / Variance Analysis, Kernel Methods

Professor Ameet Talwalkar

Slide Credit: Professor Fei Sha

Outline

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

Announcements

- HW3 and HW4 due now
- HW1 and HW2 can be picked up in class; grades available online
- Project proposal due a week from today

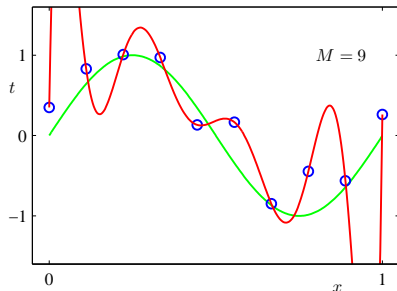
Outline

- 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

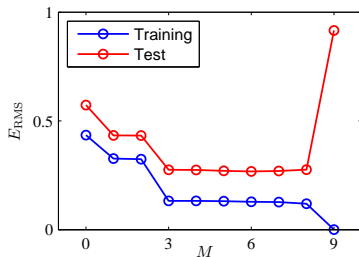
$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \Rightarrow f(x) = w_0 + \sum_{m=1}^M w_m x^m$$



Being too adaptive can improve training accuracy at the expense of test accuracy

Visualizing overfitting

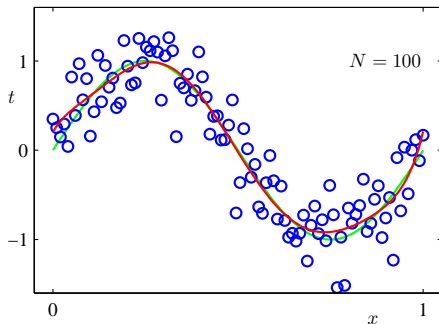
As model becomes more complex, training error keeps improving while test error first improves then deteriorates



- Horizontal axis: *measure of model complexity*, e.g., the maximum degree of the polynomial basis functions.
- Vertical axis: some notion of error, e.g., RSS for regression

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 random (*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 random (*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 random (*Bayesian* interpretation)
- To find \mathbf{w} given data \mathcal{D} , we can compute posterior distribution of \mathbf{w} :

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

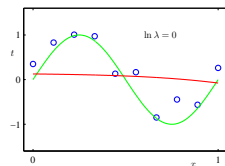
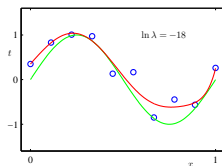
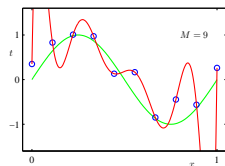
- Maximum a posterior (MAP) estimate:

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

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

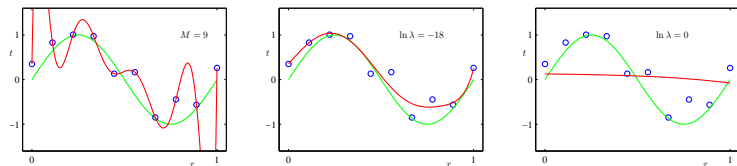
Overfitting in terms of λ

Overfitting is reduced as we increase the regularizer

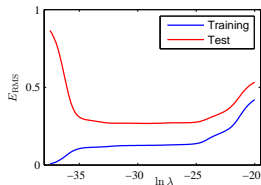


Overfitting in terms of λ

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

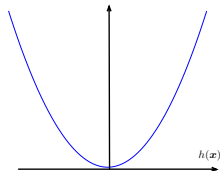
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

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

Example: quadratic loss function for regression when y is continuous

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

Ex: when $y = 0$

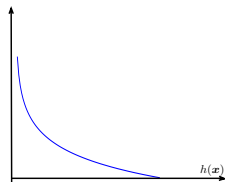


Other types of loss functions

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

$$R[h(\mathbf{x})] = \int_{\mathbf{x}, y} \ell(h(\mathbf{x}), y) p(\mathbf{x}, y) d\mathbf{x} d y$$

Measure how good our predictor is

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

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

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

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

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 this relates 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 this relates 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}
- This is called *poor generalization* or *overfitting*. We have just discussed approaches to address this issue.
- Let's try to understand 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

We are using the subscript \mathcal{D} to indicate that the prediction function is learned on the specific set of training data \mathcal{D}

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
We are using the subscript \mathcal{D} to indicate that the prediction function is learned on the specific set of training data \mathcal{D}
- $\ell(h(\mathbf{x}), y)$: our square 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
We are using the subscript \mathcal{D} to indicate that the prediction function is learned on the specific set of training data \mathcal{D}
- $\ell(h(\mathbf{x}), y)$: our square 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

Every training sample \mathcal{D} is a 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

Every training sample \mathcal{D} is a sample from the following joint distribution

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

The prediction function $h_{\mathcal{D}}(\mathbf{x})$ is a random function with respect to this distribution, and thus its risk is as well

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

The effect of finite training samples

Every training sample \mathcal{D} is a sample from the following joint distribution

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

The prediction function $h_{\mathcal{D}}(\mathbf{x})$ is a random function with respect to this distribution, and thus its risk is as well

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

Next we will disentangle the impact of the finite sample \mathcal{D} when assessing the quality of $h(\cdot)$

Average over the distribution of the training data

Averaged risk

$$\mathbb{E}_{\mathcal{D}} R[h_{\mathcal{D}}(\mathbf{x})]$$

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

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

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

Variance

We can add and subtract averaged prediction from averaged risk

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

Variance

We can add and subtract averaged prediction from averaged risk

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

Variance

We can add and subtract averaged prediction from averaged risk

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

Variance

We can add and subtract averaged prediction from averaged risk

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

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 an averaged target y

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

Bias and noise

Decompose again

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

Where is the cross-term?

Left as the take-home exercise

Analyzing the noise

How can we reduce noise

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

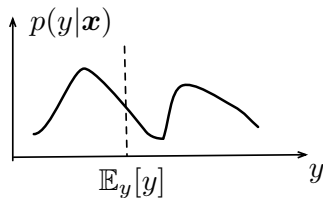
Analyzing the noise

How can we reduce noise

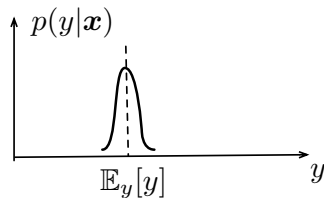
$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left(\int_y [\mathbb{E}_y[y] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

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

More difficult



Easier



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?

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

It can be reduced by using more complex models. We can choose $h(\cdot)$ to be as flexible as possible to better $h(\cdot)$ approximate $\mathbb{E}_y[y]$ and reduce bias. However, this increased flexibility will increase the VARIANCE term (as we can potentially overfit to training 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).

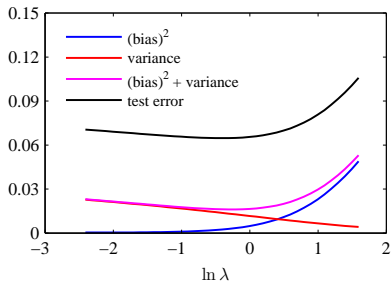
Bias/variance tradeoff

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?

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

where $\phi(\cdot)$ maps the original feature vector \mathbf{x} to a M -dimensional *new* feature vector. In the following, we will show that we can sidestep the issue of choosing which $\phi(\cdot)$ to use — instead, we will choose *equivalently* a *kernel function*.

Regularized least square

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

Motivation

How to choose nonlinear basis function for regression?

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

where $\phi(\cdot)$ maps the original feature vector \mathbf{x} to a M -dimensional *new* feature vector. In the following, we will show that we can sidestep the issue of choosing which $\phi(\cdot)$ to use — instead, we will choose *equivalently* a *kernel function*.

Regularized least square

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

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

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

where we have designated $\frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$ as α_n . And the matrix Φ is the *design matrix* made of *transformed* features. Its transpose is made of column vectors and is given by

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

where M is the dimensionality of $\phi(\mathbf{x})$.

MAP Solution

The optimal parameter vector is a linear combination of features

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

where we have designated $\frac{1}{\lambda} (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))$ as α_n . And the matrix Φ is the *design matrix* made of *transformed* features. Its transpose is made of column vectors and is given by

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

where M is the dimensionality of $\phi(\mathbf{x})$.

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

Dual formulation

We substitute $\mathbf{w}^{\text{MAP}} = \Phi^T \boldsymbol{\alpha}$ into $J(\mathbf{w})$, and obtain the following function of $\boldsymbol{\alpha}$

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

Dual formulation

We substitute $\mathbf{w}^{\text{MAP}} = \Phi^T \alpha$ into $J(\mathbf{w})$, and obtain the following 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. We see repeated structures $\Phi \Phi^T$, to which we refer as *Gram matrix* or *kernel matrix*

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

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 $\mathbf{C} = \Phi^T \Phi$
 - ① \mathbf{C} has a size of $D \times D$ while \mathbf{K} is $N \times N$.

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

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

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

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

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

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

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

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

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

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \sum_n (y - \mathbf{w}^T \boldsymbol{\phi}(x_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\Phi}^T \boldsymbol{\alpha}\|_2^2 \\ &= \frac{1}{2} \|\mathbf{y} - \mathbf{K} \boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \boldsymbol{\alpha}^T \boldsymbol{\Phi} \boldsymbol{\Phi}^T \boldsymbol{\alpha} \\ &\propto \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K}^T \mathbf{K} \boldsymbol{\alpha} - \mathbf{y}^T \mathbf{K} \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} \\ &= \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{K}^2 \boldsymbol{\alpha} - (\mathbf{K} \mathbf{y})^T \boldsymbol{\alpha} + \frac{\lambda}{2} \boldsymbol{\alpha}^T \mathbf{K} \boldsymbol{\alpha} = J(\boldsymbol{\alpha}) \end{aligned}$$

where we have used the property that \mathbf{K} is symmetric.

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)

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

Note that we only need to know \mathbf{K} in order to compute α — the exact form of $\phi(\cdot)$ is not essential — as long as we know how to get inner products $\phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$. That observation will give rise to the use of *kernel function*.

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)

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

Note that we only need to know \mathbf{K} in order to compute α — the exact form of $\phi(\cdot)$ is not essential — as long as we know how to get inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$. That observation will give rise to the use of *kernel function*.

Note that computing the parameter vector does require knowledge of Φ

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

Computing prediction needs only inner products too!

Given a test point $\phi(\mathbf{x})$, we must compute:

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

Computing prediction needs only inner products too!

Given a test point $\phi(\mathbf{x})$, we must compute:

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

where we have used the property that $(\mathbf{K} + \lambda I)^{-1}$ is symmetric (as \mathbf{K} is) and use \mathbf{k}_x as a shorthand notation for the column vector.

Note that, to make a prediction, once again, we *only need to know how to get $\phi(\mathbf{x}_n)^\top \phi(\mathbf{x})$* .

Inner products between features

Due to their central roles, let us examine more closely the inner products $\phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions consider the following $\phi(\mathbf{x})$:

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

Inner products between features

Due to their central roles, let us examine more closely the inner products $\phi(\mathbf{x}_m)^\top \phi(\mathbf{x}_n)$ for a pair of data points \mathbf{x}_m and \mathbf{x}_n .

Polynomial-based nonlinear basis functions 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,

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

Namely, the inner product can be computed by a function $(\mathbf{x}_m^\top \mathbf{x}_n)^2$ 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

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

This matrix cannot be positive semidefinite as it has both *negative* and positive eigenvalues (the sum of the diagonal elements is called the trace of a matrix, which equals to the sum of the matrix's eigenvalues. In our case, the trace is zero.)

Flashback: why using kernel functions?

without specifying $\phi(\cdot)$, the kernel 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 exactly the same as

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

‘Kernel trick’

Many learning methods depend on computing *inner products* between features — we have seen the example of regularized least squares. For those methods, we can use a kernel function in the place of the inner products, i.e., “*kernelizing*” the methods, thus, introducing nonlinear features.

We will present one more to illustrate this “trick” by kernelizing the nearest neighbor classifier.

When we talk about support vector machines, we will see the trick one more time.

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

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

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

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

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

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

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

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

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

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

The distance is equivalent to compute the distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$

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

where the $\phi(\cdot)$ is the nonlinear mapping function implied by the kernel function.

Kernelized nearest neighbors

In nearest neighbor classification, the most important quantity to compute is the (squared) distance between two data points \mathbf{x}_m and \mathbf{x}_n

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

We replace all the inner products in the distance with a kernel function $k(\cdot, \cdot)$, arriving at the kernel distance

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

The distance is equivalent to compute the distance between $\phi(\mathbf{x}_m)$ and $\phi(\mathbf{x}_n)$

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

where the $\phi(\cdot)$ is the nonlinear mapping function implied by the kernel function. The nearest neighbor of a point \mathbf{x} is thus found with

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

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.