

CSCI567 Machine Learning (Fall 2016)

Dr. Yan Liu

yanliu.cs@usc.edu

October 3, 2016

Outline

1 Bias/Variance Analysis

2 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

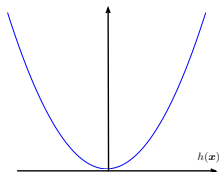
We aim to build a function $h(\mathbf{x})$ to predict the true value y associated with \mathbf{x} . If we make a mistake, we incur a *loss*

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

Example: quadratic loss function for regression when y is continuous

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

Ex: when $y = 0$

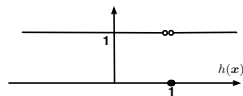


Other types of loss functions

For classification: 0/1 loss

$$\ell(h(\mathbf{x}), y) = \mathbb{I}[h(\mathbf{x}) \neq y]$$

Ex: when $y = 1$

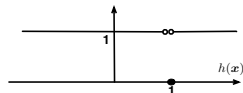


Other types of loss functions

For classification: 0/1 loss

$$\ell(h(\mathbf{x}), y) = \mathbb{I}[h(\mathbf{x}) \neq y]$$

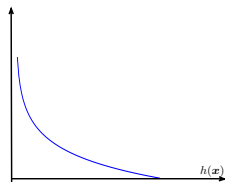
Ex: when $y = 1$



a better choice: 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

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

Measure how good our predictor is

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

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

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

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

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 learnt?

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 learnt?

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 seen how to fix that problem.

Bias/variance tradeoff

Goal: to understand the sources of prediction errors

- \mathcal{D} : our training data

Bias/variance tradeoff

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

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

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 $p(\mathbf{x}, y)$

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).

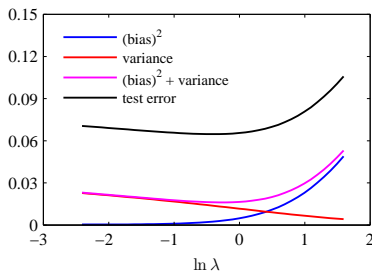
Bias/variance tradeoff

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



The effect of finite training samples

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

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

The effect of finite training samples

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

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

Thus, the prediction function $h_{\mathcal{D}}(\mathbf{x})$ is a random function with respect to this distribution. So is also its risk

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

The effect of finite training samples

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

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

Thus, the prediction function $h_{\mathcal{D}}(\mathbf{x})$ is a random function with respect to this distribution. So is also its risk

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

The first action we will take is to disentangle assessing $h(\cdot)$ from being affected by the finite samples.

Average over the distribution of the training data

Averaged risk

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

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 the predict functions learned on each training dataset.

Variance

We will subtract the averaged prediction from the 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 will subtract the averaged prediction from the 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 will subtract the averaged prediction from the 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 will subtract the averaged prediction from the 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.

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

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

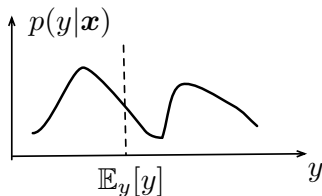
Analyzing the noise

How can we reduce noise

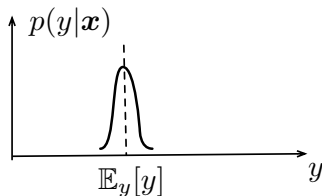
$$\int_{\mathbf{x}} \int_y [\mathbb{E}_y[y] - y]^2 p(\mathbf{x}, y) d\mathbf{x} dy = \int_{\mathbf{x}} \left(\int_y [\mathbb{E}_y[y] - y]^2 p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x}$$

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 posterior distribution $p(y|\mathbf{x})$ at the given \mathbf{x} .

Somewhat difficult posterior



Somewhat easy posterior



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 shall choose $h(\cdot)$ to be as flexible as possible: the better $h(\cdot)$ approximates $\mathbb{E}_y[y]$, the smaller the bias. However, this will increase the VARIANCE term.

Example: why regularized linear regression could be helpful?

Model

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

Consider the best possible $h^*(\mathbf{x})$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \int_{\mathbf{x}} [\mathbb{E}_y[y] - \mathbf{w}^T \mathbf{x}]^2 p(\mathbf{x}) d\mathbf{x}$$

Note that this linear model assumes the knowledge of joint distribution, thus, not achievable. Intuitively, it is the *best* linear model that can predict the data most accurately.

More refined decomposition of the bias

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

- Model bias: the price we pay for choosing linear functions to model data
- Estimation bias: the difference between the optimal model and the estimated model

Normally, the estimation bias is zero if we do not regularize.

Bias/variance tradeoff for regularized linear regression

We can only adjust estimation bias

$$\int_{\mathbf{x}} [\mathbb{E}_{\mathcal{D}} h_{\mathcal{D}}(\mathbf{x}; \lambda) - h^*(\mathbf{x})]^2 p(\mathbf{x}) d\mathbf{x}$$

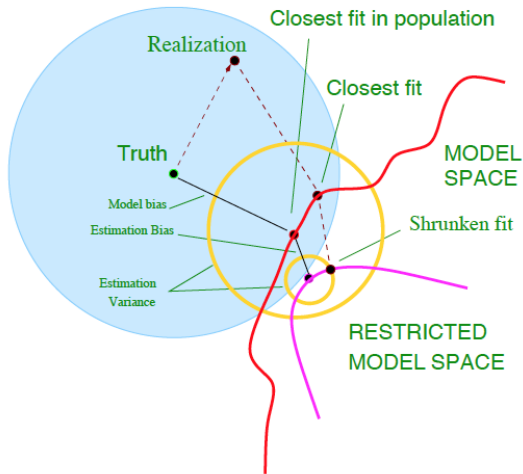
where $h(\mathbf{x}; \lambda)$ is the estimated model with regularized linear regression (parameterized with λ).

This term will not be zero anymore!

Thus, bias goes up

But, as long as this is balanced with a decrease in variance, we are willing to do so.

Visualizing the tradeoff



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).

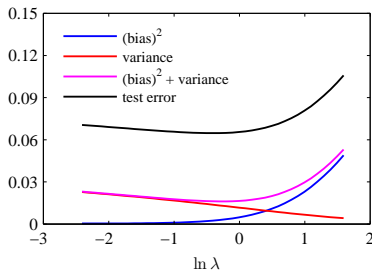
Bias/variance tradeoff

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 Bias/Variance Analysis

2 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

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

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 α (the vector of all α_n) corresponds to \mathbf{w}^{MAP} !

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*

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

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}

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

- 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 $N \times N$.
 - 2 When $N \leq D$, using \mathbf{K} is more computationally advantageous

The derivation of $J(\alpha)$

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

where we have used the property that 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)

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

Optimal α

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

which leads to (assuming that K is invertible)

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

From this, we can compute the parameter vector

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

Note that we only need to know 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 $\phi(x_m)^T \phi(x_n)$. That observation will give rise to the use of *kernel function*.

Computing prediction needs only inner products too!

Suppose we need to make a prediction on a new data point x , we thus compute

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

Computing prediction needs only inner products too!

Suppose we need to make a prediction on a new data point x , we thus compute

$$\begin{aligned} w^T \phi(x) &= \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \Phi \phi(x) \\ &= \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \begin{pmatrix} \phi(x_1)^T \phi(x) \\ \phi(x_2)^T \phi(x) \\ \vdots \\ \phi(x_N)^T \phi(x) \end{pmatrix} = \mathbf{y}^T (\mathbf{K} + \lambda I)^{-1} \mathbf{k}_x \end{aligned}$$

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 $\phi(x_n)^T \phi(x)$* .

Inner products between features

Due to their central roles, let us examine more closely the inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ 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 $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ 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,

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

Namely, the 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 $\phi(\cdot)$* .

A more challenging example

Consider the following mapping, which is parameterized by a parameter

$$\psi_{\theta}(\mathbf{x}) = \begin{pmatrix} \cos(\theta x_1) \\ \sin(\theta x_1) \\ \cos(\theta x_2) \\ \sin(\theta x_2) \end{pmatrix}$$

The inner product for transformed \mathbf{x}_m and \mathbf{x}_n is thus,

$$\psi_{\theta}(\mathbf{x}_m)^T \psi_{\theta}(\mathbf{x}_n) = \cos(\theta(x_{m1} - x_{n1})) + \cos(\theta(x_{m2} - x_{n2}))$$

Note that, once again, the inner product can be computed alternatively with the function in the right-hand-side, which depends on the original features only. We will make it further more interesting.

Concatenating into a long feature vector

Now, consider $(L + 1)$ θ s, drawn from $[0, 2\pi]$ evenly, and make another mapping

$$\phi_L(\mathbf{x}) = \begin{pmatrix} \psi_0(\mathbf{x}) \\ \psi_{\frac{2\pi}{L}}(\mathbf{x}) \\ \psi_{2\frac{2\pi}{L}}(\mathbf{x}) \\ \vdots \\ \psi_{L\frac{2\pi}{L}}(\mathbf{x}) \end{pmatrix}$$

Concatenating into a long feature vector

Now, consider $(L + 1)$ θ s, drawn from $[0 \ 2\pi]$ evenly, and make another mapping

$$\phi_L(\mathbf{x}) = \begin{pmatrix} \psi_0(\mathbf{x}) \\ \psi_{\frac{2\pi}{L}}(\mathbf{x}) \\ \psi_{2\frac{2\pi}{L}}(\mathbf{x}) \\ \vdots \\ \psi_{L\frac{2\pi}{L}}(\mathbf{x}) \end{pmatrix}$$

What is the inner product?

$$\begin{aligned} \phi_L(\mathbf{x}_m)^T \phi_L(\mathbf{x}_n) &= \sum_{l=0}^L \psi_{l\frac{2\pi}{L}}(\mathbf{x}_m)^T \psi_{l\frac{2\pi}{L}}(\mathbf{x}_n) \\ &= \sum_{l=0}^L \cos(l\frac{2\pi}{L}(x_{m1} - x_{n1})) + \cos(l\frac{2\pi}{L}(x_{m2} - x_{n2})) \end{aligned}$$

What if $L = +\infty$?

Instead of summing up $(L + 1)$ terms, we will be integrating

$$\begin{aligned}
 \phi_{\infty}(\mathbf{x}_m)^T \phi_{\infty}(\mathbf{x}_n) &= \lim_{L \rightarrow +\infty} \phi_L(\mathbf{x}_m)^T \phi_L(\mathbf{x}_n) \\
 &= \int_0^{2\pi} \cos(\theta(x_{m1} - x_{n1})) + \cos(\theta(x_{m2} - x_{n2})) d\theta \\
 &= 1 - \frac{\sin(2\pi(x_{m1} - x_{n1}))}{x_{m1} - x_{n1}} + 1 - \frac{\sin(2\pi(x_{m2} - x_{n2}))}{x_{m2} - x_{n2}}
 \end{aligned}$$

While as before, the right-hand-side depends on only the original features. It actually computes the inner product of two *infinite-dimensional* feature vectors! (Since $L \rightarrow +\infty$, the number of $\psi_{l\frac{2\pi}{L}}(\mathbf{x})$ is infinite, hence the dimensionality of $\phi_{\infty}(\mathbf{x})$.)

In other words, while we cannot write down every dimension of $\phi_{\infty}(\mathbf{x})$, we can compute its inner product easily using a function defined on the original finite feature space.

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

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

$$k(\mathbf{x}_m, \mathbf{x}_n) = 2 - \frac{\sin(2\pi(x_{m1} - x_{n1}))}{x_{m1} - x_{n1}} - \frac{\sin(2\pi(x_{m2} - x_{n2}))}{x_{m2} - x_{n2}}$$

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

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

$$k(\mathbf{x}_m, \mathbf{x}_n) = 2 - \frac{\sin(2\pi(x_{m1} - x_{n1}))}{x_{m1} - x_{n1}} - \frac{\sin(2\pi(x_{m2} - x_{n2}))}{x_{m2} - x_{n2}}$$

Examples that are not kernels

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

are not our desired kernel function as it cannot be written as inner products between two vectors.

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.

Flashback: why using kernel functions?

without specifying $\phi(\cdot)$, the kernel matrix

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

is exactly the same as

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

Examples of 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.

Examples of 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

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

Examples of 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

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

Most of those kernels have parameters to be tuned: d , c , σ^2 , etc. They are hyper parameters and are often tuned on holdout data or with cross-validation.

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.)

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, using which kernel, or which kernels to compose a new kernel, remains somewhat as “black art”, though most people will start with polynomial and Gaussian RBF kernels.

Kernelization 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., “*kernerlizing*” the methods, thus, introducing nonlinear features/basis.

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

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

Kernelized nearest neighbor classifier

In nearest neighbor classifier, 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 neighbor classifier

In nearest neighbor classifier, 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 kernelized 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 neighbor classifier

In nearest neighbor classifier, 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 kernelized 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 neighbor classifier

In nearest neighbor classifier, 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 kernelized 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

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

Take-home exercise

You have seen examples of kernelizing

- linear regression
- nearest neighbor

But can you kernelize the following?

- Decision tree
- Logistic (or multinomial logistic) regression